



UNIVERSITÀ DI PISA

---

Dipartimento di Informatica  
Corso di Laurea Triennale in Informatica

Tesi di Laurea

THE MINIMAL REPRESENTATIVE  
SELECTION PROBLEM:  
APPLICATIONS, OPTIMAL ALGORITHMS  
AND EFFICIENT HEURISTICS

Relatore:  
Prof. ANNA BERNASCONI

Candidato:  
DAVIDE RUCCI

---

ANNO ACCADEMICO 2016 - 2017



# Table of Contents

- Introduction** **iv**
  
- 1 Preliminaries** **1**
  - 1.0.1 Notation . . . . . 1
  - 1.1 Related Work . . . . . 1
  - 1.2 Characterization of optimal solutions . . . . . 3
  
- 2 Two optimal algorithms** **7**
  - 2.1 Definitions and Properties . . . . . 7
  - 2.2  $\mathcal{A}_{SM1}$ : an  $\mathcal{O}(mk)$  optimal algorithm . . . . . 11
    - 2.2.1 Correctness and complexity analysis . . . . . 14
  - 2.3  $\mathcal{A}_{SM2}$ : a faster-in-practice algorithm . . . . . 15
    - 2.3.1 Complexity Analysis . . . . . 16
  
- 3 Heuristics** **17**
  - 3.1 Pure Greedy Heuristic . . . . . 18
  - 3.2 Forward Greedy Heuristic . . . . . 19
  - 3.3 Pagerank inspired Heuristic . . . . . 20
  
- 4 Experimental Results** **25**
  
- 5 Conclusions and Future Work** **29**
  
- A Code Listing** **35**



# Introduction

In recent years, the use of *nanoscale technologies* to create electronic devices has revived interest in the use of regular structures for defining complex logic functions. One such structure is the *switching lattice*, that was initially proposed by Akers in 1972 [1], which is a wire mesh that results from laying  $m$  parallel wires perpendicularly on top of  $n$  parallel wires, with some kind of switches inserted wherever two wires cross. The kind of switch we consider is called “four-terminal” switch. A four-terminal switch is, ideally, a simple device that is linked with four neighbors and that can be either *on* or *off*: in the first case, the current arriving in the switch can flow to all of its four neighbors while in the other case, its neighboring cells are all disconnected.

These lattices can be used to implement a Boolean function by, firstly, connecting all the wire ends at the top to a “top plate”  $T$  and all the wire ends at the bottom to a “bottom plate”  $B$  and, secondly, having each four-terminal switch controlled by an input. An input wire can be inverted, thus allowing to associate each four-terminal switch to a Boolean literal or its negation, namely  $x$  or  $\neg x$ . Some input combinations will create a (possibly meandering) path between top and bottom plates: we associate those combinations with output being 1, or *true*, while every combination that leaves  $T$  and  $B$  unconnected is associated with the output 0, or *false*. In a similar way, we can also equip a switching lattice with “left” and “right” plates,  $L$  and  $R$  respectively, and the left-to-right connectivity can be associated to a Boolean function as well. In this way, a single lattice can implement two different Boolean functions, although not simultaneously. Altun and Riedel, in 2012, exploited this possibility to develop a new synthesis algorithm [2] for this kind of switching lattices, that is, an algorithm that, given a target Boolean function, assigns a literal to each four-terminal switch in a way such that there is connectivity between  $T$  and  $B$  if and only if the target function evaluates to true.

If we picture the  $m$  wires laid down horizontally, we can visualize a matrix made by  $m$  rows and  $n$  columns of switches, as shown in Figure 1.

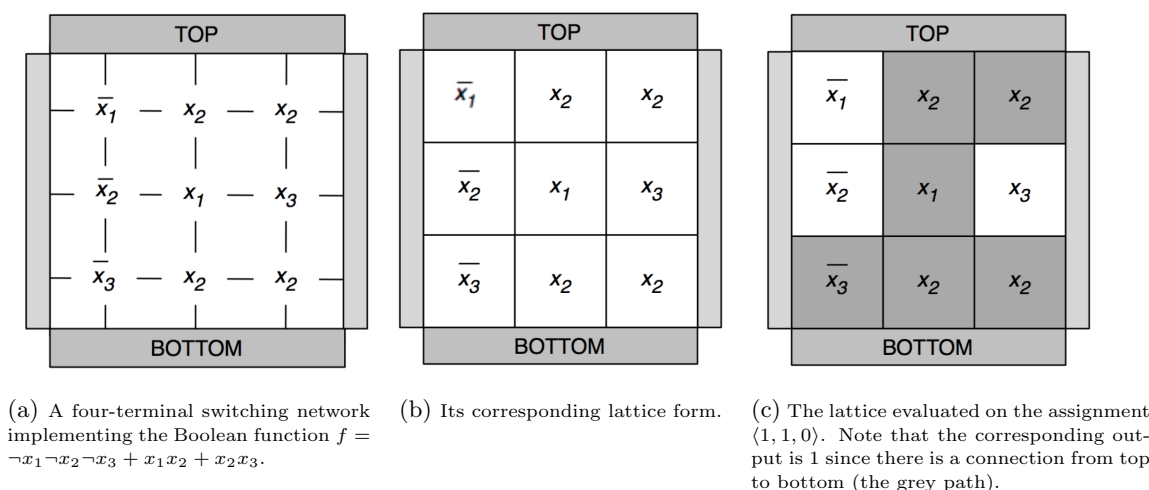


Figure 1: An example of a switching lattice.

We now describe the algorithm proposed in [2] for the synthesis of a  $m \times n$  switching lattice. Let  $f_T$  be the function to be implemented by the lattice, and let  $f_T^D$  be the *dual function* of  $f_T$ , i.e., the function obtained from  $f_T$  by interchanging the *AND* and *OR* operations, as well as interchanging the 0 and 1

constants. All the functions are represented in ISOP form, which stands for *Irredundant Sum-Of-Products*. Informally, a Boolean function  $f$  is in ISOP form if it is expressed by an *OR* (sum) of *ANDs* (products) and if no product can be removed from it without changing  $f$ . For example,  $f = (x_1x_2) + (\neg x_1x_3)$  is in ISOP form, and its dual function is  $f^D = (x_1 + x_2)(\neg x_1 + x_3)$ .

The algorithm is as follows (suppose that  $f_T$  and  $f_T^D$  have  $n$  and  $m$  products, respectively):

- Start with a  $m \times n$  lattice;
- Assign each product of  $f_T$  to a column, and each product of  $f_T^D$  to a row;
- Compute the intersection<sup>1</sup> of the literals in the corresponding products.
- From the resulting sets associated with each cell, pick one literal arbitrarily and assign it to the corresponding site of the lattice.

This way, there will be a connected path from  $T$  to  $B$  if and only if  $f_T$  evaluates to 1. A proof of the correctness of this algorithm can be found in [2, Section 3.2], while Figure 2 provides a visualization of it.

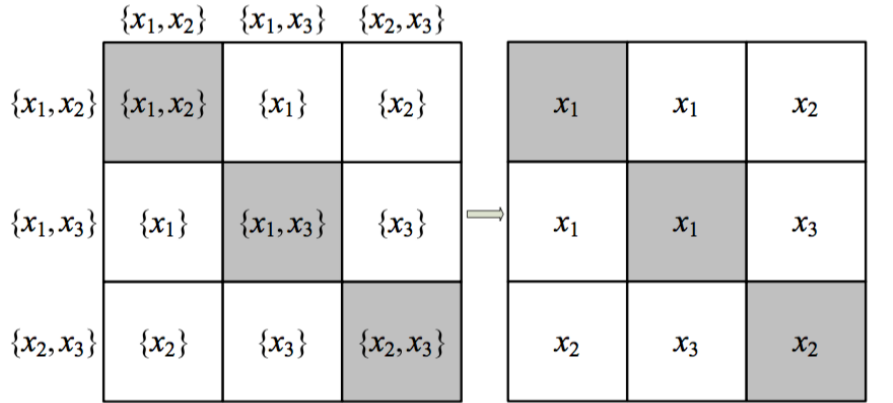


Figure 2: Left: the intersections computed by the algorithm for the function  $f_T = x_1x_2 + x_1x_3 + x_2x_3$  on a  $3 \times 3$  switching lattice. Note that  $f_T^D = f_T$ . Right: an arbitrary assignment of the variables to their corresponding sites on the lattice.

Our work is motivated by the last step of the algorithm, specifically by the way in which the literals are chosen from the sets created after doing the rows-columns intersection. When synthesizing this kind of lattices it may be a good idea to not pick too many times the same literal, because plugging the same literal in many different cells has a higher cost, in terms of current supply needed, than plugging, say, many different literals in fewer cells. For this reason, our work focuses on *minimizing the maximum frequency* of all the literals, which we will call *variables*, across all the sets.

Although this problem has not been directly studied in the past, we have found some interesting connections to many different, and known, problems such as graphs orientation problems, networks flow problems and scheduling problems. In fact, the latter area gave us the most important results from a theoretical point of view, since it has allowed us to state that our problem can be solved in polynomial time by providing two algorithms that are able to solve always correctly our problem. However, we have put the majority of our efforts in developing heuristics for this problem, for basically two good reasons. The first focuses on execution times: heuristics, by their nature, tend to be way faster than “standard” algorithms at the cost of possibly introducing errors in the solutions they give; the heuristics that we propose, and in particular the first two, are very fast even with large input instances and they have shown, in our experiments, a success rate high enough to induce one to prefer them instead of an always-optimal algorithm. The second reason, which is in some sense the most important within this context, is that one of the two optimal algorithms that we are going to describe in Chapter 2 bases his work on an initialization with a heuristic. In other words, this algorithm takes an arbitrary, possibly non-optimal, solution and repeatedly improves it until it becomes optimal. It is clear that the better its initialization, the less time the algorithm will take to produce an optimal solution.

<sup>1</sup>The intersection will always contain at least one element, as shown in [2, Lemma 1].

# Chapter 1

## Preliminaries

Let  $\mathcal{N} = \{x_1, x_2, \dots, x_n\}$  be a set of  $n$  variables and let  $\mathcal{F}$  be a family of  $m$  subsets of  $\mathcal{N}$ . We consider the problem of finding a representative for each subset in  $\mathcal{F}$  such that the frequency of the most chosen one is minimal, i.e., the maximum frequency is minimal. In general,  $\mathcal{F}$  can contain repeated sets. We call this problem *Minimal Representative Selection*, or **MRS**, which can be formally defined as follows:

**Definition 1.1** (MRS). Given  $\mathcal{N}$  and  $\mathcal{F}$ , the **MRS** problem is:

INSTANCE: a set  $\mathcal{N}$  of  $n$  variables and a family  $\mathcal{F}$  of  $m$  subsets of  $\mathcal{N}$ .

QUESTION: find an assignment of representatives for every subset in  $\mathcal{F}$  such that the frequency of the most selected one is minimized.

We will show that **MRS** is in  $\mathcal{P}$  by describing two algorithms [4] that solve it in polynomial time. These two algorithms come from the area of scheduling problems. In fact, we will show how this problem can be reduced to, at least, three different problems, the first one is valid only for a restricted version of **MRS** [5], while the other two can solve every instance of it [4, 6].

For various reasons, we can put weights on the subsets in  $\mathcal{F}$ , but this will make **MRS**  $\mathcal{NP}$ -hard (see the the reduction from PARTITION in [5]).

### 1.0.1 Notation

The size of an input instance of **MRS** is defined as  $k = \sum_{i=1}^m |\mathcal{F}[i]|$  where  $\mathcal{F}[i]$  is the  $i$ th set in  $\mathcal{F}$ . We will refer to the *value*<sup>1</sup> of the optimal solution as  $t^*$ , while the value of a generic solution will be denoted by  $t$ . Let  $S$  be a subset of  $\mathcal{N}$  and let  $\mathcal{F}_{|S} = \{A \in \mathcal{F} \mid \forall a \in A \implies a \in S\}$ . We call  $\mathcal{F}_{|S}$  the *subfamily induced by  $S$* . To lighten the notation, we denote the cardinality of the two sets  $\mathcal{F}_{|S}$  and  $S$  as  $m_S$  and  $n_S$  respectively.

## 1.1 Related Work

The **MRS** problem as defined above has not been directly studied in the past. However, there are some interesting points in common with problems from different areas such as graph orientation problems [5], network flow problems [6] and a variant of bipartite maximum matching called *semi-matching*, used to solve a particular type of the *offline scheduling problem* [4]. These problems are all solvable in polynomial time, making our problem polynomially solvable too, thanks to three simple reductions.

We first explain the reduction to a graph orientation problem that is possible, however, only for a restricted version of **MRS**, which we call **2-MRS**. In this restricted version, we ask that all subsets in  $\mathcal{F}$  have cardinality *exactly* 2 and are distinct. The instances of this problem can be then modeled with a graph as follows. Let  $G = (V, E)$  be a simple, undirected, weighted graph with  $V = \mathcal{N}$  and  $E = \mathcal{F}$ . We can say that choosing a variable from each subset, represented by the edges in  $G$ , is equivalent to

---

<sup>1</sup>The value of a solution of **MRS** is the frequency of the most chosen variable as representative.

choosing an orientation for every member of  $E$ . Since we want to minimize the maximum frequency of the representatives, we have to minimize the maximum indegree or outdegree of the vertices of  $G$  depending on how we choose to orient the edges. We focus on minimizing the *outdegree*, meaning that the orientation of the edges will be from the variable chosen as representative of each subset, to the other variable paired with it. The results we report here come from [5], while the corresponding results for the indegree version can be found in [7].

The *Minimum Maximum Outdegree Orientation* problem, or *MMO* is, formally, the problem of giving an orientation to every edge of a weighted graph which minimizes the maximum weighted outdegree of its vertices and it is in general  $\mathcal{NP}$ -hard, but in the special case of all equal weights, the problem becomes polynomial. Since we have not put weights on the subsets in  $\mathcal{F}$ , we can consider the weight function of  $G$  as  $w(\{i, j\}) = 1$  for every  $\{i, j\} \in E$  and we can solve **2-MRS** in  $\mathcal{O}(m^2)$  time thanks to the algorithm REVERSE described in [5]. The basic strategy of the aforementioned algorithm is to find first an arbitrary orientation of edges and then repeatedly improve it by reversing the orientation of edges along a path that satisfies two inequalities related to the outdegree of its first and last vertices. Figure 1.1 shows an example of the graph corresponding to the **2-MRS** instance defined by  $\langle \mathcal{N} = \{1, 2, 3, 4, 5\}, \mathcal{F} = \{\{1, 2\}, \{1, 3\}, \{2, 3\}, \{2, 4\}, \{4, 5\}\} \rangle$  and its optimal solution (1, 3, 2, 4, 5).

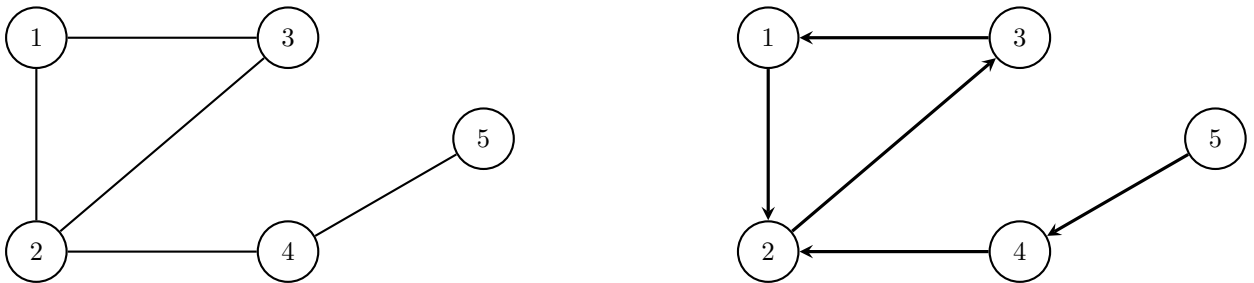


Figure 1.1: The graph of an instance of **MRS** (left) and the solution obtained by applying the REVERSE algorithm to it (right).

The area of scheduling problems is, however, the one that provides us with the most useful reductions for solving every instance of **MRS**. The scheduling problem that is equivalent to **MRS** is a restricted version of the *offline Parallel Machine Scheduling* problem, or *PMS*, which has received intensive study in the literature (see, for example, [8,9]). The parallel machine scheduling problem, denoted by the so-called *three-field representation*  $R \mid \circ \mid C_{\max}$ <sup>2</sup> [10], is the problem of assigning  $m$  jobs to  $n$  identical, unrelated, machines. The goal is to minimize one or more objectives such as the *makespan*<sup>3</sup> of the schedule or the *flow time*<sup>4</sup> of the jobs. If all the jobs require one unit of time to complete and can be executed only on a specific subset of the available machines, the problem is denoted by  $R \mid p_j = 1, \mathcal{M}_j \mid C_{\max}$ <sup>5</sup>. Minimizing the makespan is, in particular, the main objective of **MRS**, as it is equivalent to minimize the maximum frequency of all the representatives chosen.

We first show how **MRS** can be reduced to PMS with unit-length jobs. Every variable in  $\mathcal{N}$  is a machine, and every subset in  $\mathcal{F}$  can be thought as a job that can be executed by the machines denoted by the variables it contains. For example, the subset  $\{1, 2, 3\}$  of  $\mathcal{N} = \{1, 2, 3, 4, 5\}$  represents a unit-length job that can be executed either on machine 1, 2 or 3, but not on machines 4 or 5. More formally, we can define  $J = \mathcal{F}$  as the set of jobs,  $\mathcal{M} = \mathcal{N}$  as the set of machines and fix  $p_j = 1$  for every  $1 \leq j \leq m$ . The solution of PMS, that is an assignment of jobs to machines, can be seen as an assignment of sets to a variable, thus creating the choice of representatives for each subset in  $\mathcal{F}$ .

The  $R \mid p_j = 1, \mathcal{M}_j \mid C_{\max}$  problem, obtained by the reduction from **MRS**, can be solved by a polynomial

<sup>2</sup>For the sake of clarity,  $R$  stands for *unrelated* parallel machines,  $\circ$  indicates that there are no particular restrictions on the jobs, and  $C_{\max}$  is the objective we want to reach, that is the one of minimizing the makespan.

<sup>3</sup>The *makespan* of a schedule is the difference, in time, between its start and its finish. With unit-length jobs, it represents the maximum number of jobs assigned to a machine.

<sup>4</sup>The *flow time* of a job is defined as the time at which the job is completed minus the time at which the job was made available for processing to its machine.

<sup>5</sup>Note that the second field is no longer empty: it contains a restriction on the time that it takes for each job to complete ( $p_j = 1$ ) and a restriction on the machines that are capable to execute the job  $j$  ( $\mathcal{M}_j \subseteq \mathcal{M} =$  the set of all the machines available).

time algorithm for solving the *network maximum flow* problem (see, for example, algorithms in [11]), thus involving a further reduction to graphs. Let  $G = (V, E, w)$  be a simple, directed, and capacitated graph with  $V = \mathcal{F} \cup \mathcal{N} \cup \{s, t\}$ . We denote by  $x_i$ , for  $1 \leq i \leq n$ , the vertices that represent a variable in  $\mathcal{N}$  and by  $y_j$ , for  $1 \leq j \leq m$ , the vertices that represent a subset in  $\mathcal{F}$ . Since  $s$  and  $t$  will be the source and the sink of the maximum flow, respectively,  $E$  will be constructed as follows:

- an edge  $(s, x_i)$  for every variable-vertex  $x_i$ ;
- an edge  $(x_i, y_j)$  if and only if the variable represented by vertex  $x_i$  is in the corresponding subset of  $\mathcal{F}$  represented by the vertex  $y_j$ ;
- an edge  $(y_j, t)$  for every subset-vertex  $y_j$ .

Then we define the *capacity function*  $w : E \rightarrow \mathbb{N}$  as:

- $w(s, x_i) = c$  where  $c$  is a parameter;
- $w(x_i, y_j) = w(y_j, t) = 1$  for every  $i$  and  $j$ .

The following theorem correlates the maximum flow problem and PMS. The proof can be found in [6].

**Theorem 1.1.** *The problem  $R|p_j = 1, \mathcal{M}_j|C_{max}$  has a feasible makespan of  $C_{max} \leq c$  if and only if the value of the maximum flow in  $G$  is  $m$ .*

The algorithm proposed in [6] is a simple “binary search” of the minimum possible value of  $c$ . We report the pseudocode here:

---

**Algorithm 1** (Binary Search)

---

```

1: function BS( $G$ )
2:    $l \leftarrow 1$ 
3:    $u \leftarrow m$ 
4:   repeat
5:      $c \leftarrow \lfloor \frac{l+u}{2} \rfloor$ 
6:      $f \leftarrow$  The maximum flow in  $G$ 
7:     if  $f = m$  then ▷  $C_{max}^* \leq c$ 
8:        $u \leftarrow c$ 
9:     else ▷  $C_{max}^* > c$ 
10:       $l \leftarrow c + 1$ 
11:   until ( $l = u$ )
12:   return The schedule of makespan  $l$  ▷ The optimal  $c$  value is reached ( $C_{max}^* = l$ )

```

---

**Theorem 1.2.** *The Binary Search algorithm correctly solves the  $P|p_j = 1, \mathcal{M}|C_{max}$  problem in  $\mathcal{O}(m^3 \log m)$  time. [6]*

Although this is a polynomial algorithm that correctly solves every instance of **MRS**, we are not completely satisfied by its time bound, cubic in  $m$ , which can be very large in practice and sometimes it can be even exponential in  $n$ . This is because the main strategy of this algorithm, i.e., finding the maximum flow in  $G$ , is not specialized on solving this kind of problems. For this reason, in the next chapter we describe two algorithms with lower running times:  $\mathcal{O}(mk) = \mathcal{O}(m^2n)$  and  $\mathcal{O}(m^2n^2)$  respectively. The second bound, however, is loose and we will show that in practice, the latter algorithm performs even better than the former.

## 1.2 Characterization of optimal solutions

We provide three theoretical results on **MRS**. The first two are bounds on the solutions, that can be used for developing good heuristics, and the third allows us to simplify the input instances, i.e., to reduce the size of the input. Recall that  $t$  and  $t^*$  denote the value of a generic and an optimal solution respectively.

**Proposition 1.2.1** (Weak Lower Bound).  $t^* \geq \lceil \frac{m}{n} \rceil$

*Proof.* Let  $r$  be the number of representatives chosen in a solution of **MRS**. Clearly  $r = m$ . For the sake of contradiction, suppose that  $t^* < \lceil \frac{m}{n} \rceil$ . We now distinguish two cases: first assume that  $n \mid m$ . This implies that  $t^* < \frac{m}{n}$ . We can give an upper bound on  $r$  by saying that we choose each variable in  $\mathcal{N}$  as representative at most  $t^*$  times.

$$r \leq n \cdot t^* < n \cdot \left\lceil \frac{m}{n} \right\rceil = n \cdot \frac{m}{n} = m$$

with the strict inequality coming from our hypothesis on  $t^*$ . This leads us to conclude that  $r < m$ , which is a contradiction.

Conversely, if  $n \nmid m$ , then  $t^* < \lceil \frac{m}{n} \rceil \implies t^* \leq \lfloor \frac{m}{n} \rfloor$ , so

$$r \leq n \cdot \left\lfloor \frac{m}{n} \right\rfloor < n \cdot \frac{m}{n} = m$$

which is the same contradiction as above. In other words we left at least one set without its representative, thus  $t^* \geq \lceil \frac{m}{n} \rceil$ .  $\square$

In a similar way, we can improve this lower bound by considering all the subfamilies induced by any subset  $S \subseteq \mathcal{N}$  on  $\mathcal{F}$ .

**Theorem 1.3** (Lower Bound).  $t^* \geq \max_{S \subseteq \mathcal{N}} \left\lceil \frac{m_S}{n_S} \right\rceil$

*Proof.* We distinguish two cases as before. First, assume that  $n_S \mid m_S$  and suppose, for the sake of contradiction, that there exists  $S \subseteq \mathcal{N}$  such that  $t^* < \lceil \frac{m_S}{n_S} \rceil$ . This implies that  $t^* < \frac{m_S}{n_S}$ . Now, using the same argument as in Proposition 1.2.1, we can obtain the following upper bound on the number of representatives chosen in  $\mathcal{F}_S$ :

$$m_S \leq n_S \cdot t^* < n_S \cdot \frac{m_S}{n_S} = m_S$$

which is a contradiction. Conversely, if  $n_S \nmid m_S$ , the hypothesis  $t^* < \lceil \frac{m_S}{n_S} \rceil$  implies that  $t^* \leq \lfloor \frac{m_S}{n_S} \rfloor$ . The upper bound now is

$$m_S \leq n_S \cdot t^* \leq n_S \cdot \left\lfloor \frac{m_S}{n_S} \right\rfloor < n_S \cdot \frac{m_S}{n_S} = m_S$$

which is also a contradiction, with the strict inequality coming from the well-know properties of the ceiling and floor functions. This means that not reaching this lower bound on every possible  $\mathcal{F}_S$  implies leaving at least one set without its representative.  $\square$

This theorem follows the intuition that, when there are many singletons of the same variable, the lower bound is at least the number of occurrences of that singleton. In the other cases, it allows us to say that there exists a ‘‘difficult core’’ in the initial family that inevitably increases the global lower bound. It also allows us to simplify the instances of **MRS** by immediately choosing as representatives the variables that appear less then this lower bound without altering the value of the optimal solution, as explained with the following Corollary.

**Corollary 1.3.1** (Input Simplification). *If a variable in  $\mathcal{N}$  appears less than  $t^*$  times in  $\mathcal{F}$  then there exists an optimal solution in which that variable is chosen as representative in every subset containing it.*

*Proof.* Suppose that the variable  $x$  appears in  $z < t^*$  sets of  $\mathcal{F}$ . We are going to prove that there exists an optimal solution with  $x$  being chosen exactly  $z$  times by giving an algorithm that constructs this solution from an optimal solution already obtained. Let  $W$  be an optimal solution for **MRS** with  $x$  chosen less than  $z$  times and apply the following algorithm to it:

1. Sort the variables according to the number of times they have been chosen as representative, in ascending order, ignoring the variable  $x$ ;
2. Loop through the sorted variables:

- (i) If the current variable appears in a set containing  $x$ , change the representative of this set to be  $x$ , otherwise consider the next variable;
- (ii) If  $x$  has been chosen  $z$  times, stop.

This algorithm never changes the representative of a set which contains  $x$  and the *last* most chosen variable, otherwise  $W$  would not have been an optimal solution. It also cannot worsen the value of  $W$  because it only decreases frequencies of the other variables and increases the one of  $x$  which is bounded by  $z < t^*$ .  $\square$

It is important to note that this algorithm can be recursively applied to the instance it outputs, because the instance created is actually a new, *independent*, one with different values of  $m$  and  $n$ .

Before proceeding, we give some practical considerations. This kind of simplification is generally possible with gaussian-distributed variables across the sets in  $\mathcal{F}$ , if  $m$  is sufficiently large. This is likely to happen because, given the nature of **MRS**,  $m$  can be even exponential in  $n$ , thus increasing the value of  $\lceil \frac{m}{n} \rceil$  which is, almost always, the only lower bound efficiently computable<sup>6</sup>. With a large  $\lceil \frac{m}{n} \rceil$  and a gaussian distribution of the variables, there will be at least one “removable” variable in many instances. The same reasoning cannot be applied when the variables are uniformly distributed across the sets in  $\mathcal{F}$  because, for large values of  $m$ , a variable will tend to appear the same number of times with respect to another one, making the simplification to be roughly impossible.

---

<sup>6</sup>A lower bound higher than the one given in Proposition 1.2.1 can still be obtained by counting the occurrences of the most-frequent singleton.



# Chapter 2

## Two optimal algorithms

In this chapter we describe two algorithms for computing an optimal solution of **MRS**. These algorithms have been designed to compute an optimal semi-matching in bipartite graphs [4]. We first define semi-matchings, then we describe some of their properties and the two optimal polynomial-time algorithms.

Let  $G = (U \cup V, E)$  be an undirected, bipartite graph, in which  $E \subseteq U \times V$ . A *matching*  $M$  in  $G$  is a set of edges,  $M \subseteq E$ , such that each vertex in  $U \cup V$  is an endpoint of at most one edge in  $M$ . In the following we will say that a vertex  $u \in U$  is *matched* with the vertex  $v \in V$  if there is an edge in  $M$  from  $u$  to  $v$  and vice-versa.

**Definition 2.1** (Semi-Matching). Let  $G = (U \cup V, E)$  be a bipartite graph. A semi-matching in  $G$  is a subset of edges,  $M \subseteq E$ , such that each vertex  $u \in U$  is matched with *exactly* one vertex  $v \in V$ .

**Remark.** A semi-matching cannot exist if  $G$  has at least an isolated  $U$ -vertex.

Note that finding a semi-matching in  $G$  is trivial, just match each vertex in  $U$  with one vertex of  $V$ . Recall that semi-matchings in [4] are used to solve an offline scheduling problem, in which the optimization objectives are many and different from each other. One is to minimize the *makespan* of the schedule, while another possible objective can be to minimize the average completion time, of the tasks. The following algorithms were designed to reach these many different optimization objectives at the same time and we will show how this can be achieved in practice.

The scheduling problem defined in Section 1.1 can be translated in terms of semi-matchings as follows: each task is represented by a vertex  $u \in U$  and each machine is represented by a vertex  $v \in V$ . There is an undirected edge  $\{u, v\}$  if the task  $u$  can be executed on the machine  $v$ . The edges in the semi-matching found mean to assign job  $u$  to machine  $v$ .

### 2.1 Definitions and Properties

Recall the equivalence of  $U$  and  $V$  to our  $\mathcal{F}$  and  $\mathcal{N}$ , respectively. For  $v \in V$ , let  $\deg(v)$  denote the degree of vertex  $v$  or, in load balancing terms, the number of tasks that machine  $v$  is capable to process. Given a semi-matching  $M$ , we denote with  $\deg_M(v)$  the degree of  $v$  in  $M$ , but we will frequently refer to it as the *load* on vertex  $v$ . Since each task requires one unit of time to complete, we can define the *cost* of a vertex  $v \in V$  as its jobs' *flow time*.

**Definition 2.2** (Vertex Cost).  $\text{cost}_M(v) = \sum_{i=1}^{\deg_M(v)} i = \frac{(\deg_M(v)+1)\deg_M(v)}{2}$

This is because we suppose that all jobs are available for processing to their assigned machine at time 0. In a similar way, we define the *total cost* of a semi-matching  $M$  as the sum of all  $V$ -vertices' cost.

**Definition 2.3** (Total Cost).  $C(M) = \sum_{i=1}^n \text{cost}_M(v_i)$

A semi-matching is optimal when its cost is minimized and an optimal semi-matching is also optimal with respect to other optimization objectives, as shown later in this section.

The following definitions create the basis for the theoretical results that will be used in the next algorithms.

**Definition 2.4** (Alternating Path). Let  $M$  be a semi-matching in  $G$ . An alternating path is a sequence of edges  $P = (\{v_1, u_1\}, \{u_1, v_2\}, \{v_2, u_2\}, \dots, \{u_{k-1}, v_k\})$  in which  $v_i \in V$ ,  $u_i \in U$  and  $\{v_i, u_i\} \in M$  for each  $1 \leq i \leq k$ .

For convenience, in the following we often treat paths as a sequence of vertices instead of a sequence of edges, i.e.,  $P = (v_1, u_1, v_2, u_2, \dots, u_{k-1}, v_k)$ . So far we have considered an undirected graph, but in the following we assume instead that edges in  $E$  are directed. More precisely, the edges of  $G$  in a semi-matching  $M$  will be directed from a vertex  $v \in V$  to a vertex  $u \in U$  while the others will be directed from  $U$  to  $V$ . When switching matching and non-matching edges along path  $P$  we create a new semi-matching with the load on  $v_1$  decreased by one, the load on  $v_k$  increased by one and the load of the other  $v \in V$  unchanged. This derived semi-matching is denoted by  $P \oplus M$ , because it is formally equivalent to the symmetric difference of sets, i.e.,  $A \oplus B = (A \setminus B) \cup (B \setminus A)$ .

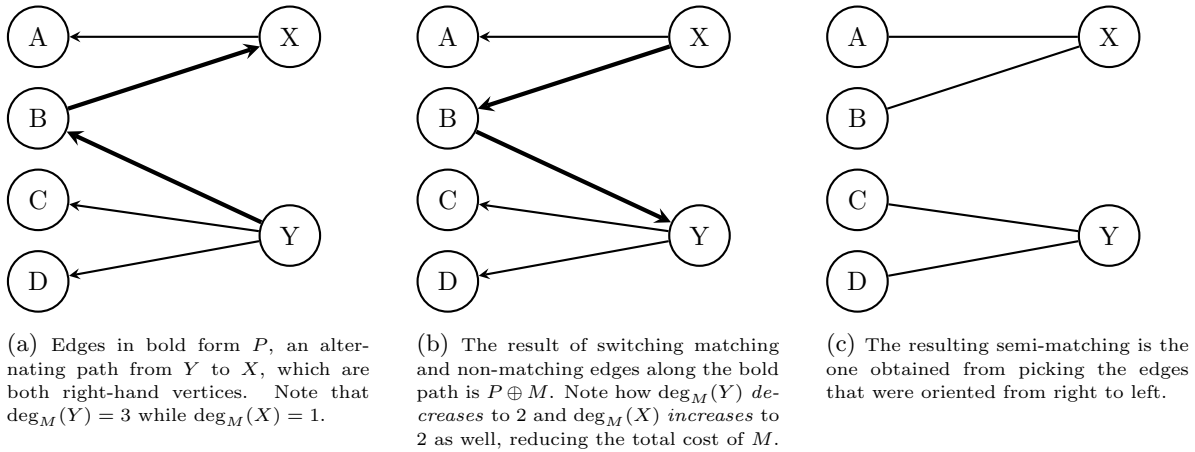


Figure 2.1: An example of alternating path.

**Definition 2.5** (Cost-reducing Path). An alternating path in which  $\deg_M(v_1) > \deg_M(v_k) + 1$  is called a *cost-reducing path*.

Cost-reducing paths are so-called because the new semi-matching  $P \oplus M$  will have cost  $C(P \oplus M) = C(M) - (\deg_M(v_1) - \deg_M(v_k) - 1)$  which is less than the previous cost,  $C(M)$ . For example, the bold path in Figure 2.1(a) is a cost-reducing path. These type of paths represent a solid theoretical base of semi-matchings as we are going to show with the following theorem [4].

**Theorem 2.1.** Let  $G = (U \cup V, E)$  be a bipartite graph. A semi-matching  $M$  in  $G$  is optimal if and only if no cost-reducing path exists in it.

*Proof.* It is clear that if  $M$  is optimal, no cost-reducing path can exist, otherwise  $M$  would not have been optimal because the load on its first vertex could be reduced by 1. Therefore we prove that a cost-reducing path must exist if  $M$  is not optimal.

Let  $O$  be an optimal semi-matching in  $G$ , chosen such that the symmetric difference  $O \oplus M$  is minimized and suppose that  $M$  is not optimal. This implies that  $M$  has greater total cost than  $O$ , i.e.,  $C(O) < C(M)$ . Let  $G_d$  be the subgraph induced by the edges of  $O \oplus M$  and color with green the edges of  $O \setminus M$  and with red the ones of  $M \setminus O$ . Direct the green edges from  $U$  to  $V$  and the red edges from  $V$  to  $U$ .

To complete the proof we first state and prove the following proposition.

**Proposition 2.1.1.** The graph  $G_d$  is acyclic and for every directed path  $P$  in  $G_d$  from  $v_1 \in V$  to  $v_2 \in V$ ,  $\deg_O(v_2) \leq \deg_O(v_1)$  holds.

*Proof.* Let  $P = (v_1, \dots, v_2)$  be a directed path in  $G_d$ . By the choice of directions for the edges,  $P$  must be made of alternating red and green edges. If  $\deg_O(v_2) > \deg_O(v_1) + 1$  then the reverse of  $P$  is a cost-reducing path for  $O$ , in contradiction to the optimality of  $O$ . Conversely, if  $\deg_O(v_2) = \deg_O(v_1) + 1$ , then we can get an optimal assignment with smaller symmetric difference from  $M$  by alternating the

assignment of the  $U$ -vertices along  $P$ , in contradiction to the minimality of  $O \oplus M$ . Therefore it must be that  $\deg_O(v_2) \leq \deg_O(v_1)$ .

A similar argument shows that  $G_d$  is also acyclic.  $\blacksquare$

Recall that both  $O$  and  $M$  are semi-matchings, implying that  $\sum_{v \in V} \deg_O(v) = \sum_{v \in V} \deg_M(v) = |U|$ . Since  $C(O) < C(M)$ , there exists  $v_1 \in V$  such that  $\deg_M(v_1) > \deg_O(v_1)$ . Starting from  $v_1$ , we can build an alternating red-green path  $P'$  as follows:

1. From an arbitrary vertex  $v \in V$ , if  $\deg_{M \setminus O}(v) \geq 1$  and  $\deg_M(v) \geq \deg_M(v_1) - 1$ , we build  $P'$  by following an arbitrary red edge directed out from  $v$ .
2. From an arbitrary vertex  $u \in U$ , we build  $P'$  by following the single green edge directed out from  $u$ . Such a green edge must always exist, because  $O$  and  $M$  are semi-matchings.
3. Otherwise, stop.

By Proposition 2.1.1,  $G_d$  is acyclic and therefore  $P'$  is well defined and finite.

Now let  $v_2 \in V$  be the final vertex on the path. We distinguish two cases: if  $\deg_M(v_2) < \deg_M(v_1) - 1$  then  $P'$  is a cost-reducing path relative to  $M$ . If  $\deg_{M \setminus O}(v_2) = 0$ , we know that  $\deg_M(v_2) < \deg_O(v_2)$  since  $P'$  arrived at  $v_2$  by following a green edge. By Proposition 2.1.1 we also know that  $\deg_O(v_2) \leq \deg_O(v_1)$ . Recall that  $v_1$  was chosen such that  $\deg_O(v_1) < \deg_M(v_1)$ . Combining these three inequalities yields  $\deg_M(v_2) < \deg_O(v_2) \leq \deg_O(v_1) < \deg_M(v_1)$ , implying that  $\deg_M(v_2) < \deg_M(v_1) - 1$ , hence  $P'$  is a cost-reducing path relative to  $M$ .

Since  $P'$  is a cost-reducing path relative to  $M$  in both cases, the proof is complete.  $\square$

So far we have considered the objective of minimizing only  $C(M) = \sum_{i=1}^n \text{cost}_M(v_i)$ , where  $\text{cost}_M(v) = \sum_{i=1}^{\deg_M(v)} i$  which is, in other words, the flow time of the tasks' schedule. We now prove that we can extend these results to a wider set of cost functions. First, we formally define a *cost function*.

**Definition 2.6** (Cost Function). A cost function for a matching  $M$  is a function

$$\text{cost}_f(M) = \sum_{i=1}^n f(\deg_M(v_i))$$

where  $f : \mathbb{R}^+ \rightarrow \mathbb{R}$  is a *strictly convex* function. If  $f$  is weakly convex, then  $\text{cost}_f(M)$  is called a *weak cost function*.

For example, our  $C(M)$  is the cost function induced by  $f(x) = \sum_{j=1}^x j = \frac{x(x+1)}{2}$ .

We can now extend the concept of cost-reducing paths according to this definition.

**Lemma 2.1.** *Let  $P$  be a cost-reducing path. For any cost function  $\text{cost}_f$ , switching matching and non-matching edges along  $P$  yields a semi-matching  $M' = P \oplus M$  such that  $\text{cost}_f(M') < \text{cost}_f(M)$ . If  $\text{cost}_f$  is a weak cost function then  $\text{cost}_f(M') \leq \text{cost}_f(M)$ .*

*Proof.* First, we recall the following property of convex functions:

**Proposition 2.1.2.** *If  $f$  is strictly convex then  $x > y \implies f(x+1) - f(x) > f(y+1) - f(y)$ . If  $f$  is weakly convex then the inequality holds but is not strict.*

*Proof.* A property of strictly convex functions is that

$$\frac{f(b) - f(a)}{b - a} < \frac{f(c) - f(a)}{c - a} < \frac{f(c) - f(b')}{c - b'}$$

if  $a < b < c$  and  $a < b' < c$  [12]. Setting  $a = y, b = y + 1, b' = x$  and  $c = x + 1$  we obtain

$$f(y+1) - f(y) < \frac{f(x+1) - f(y)}{x+1-y} < f(x+1) - f(x)$$

If  $f$  is weakly convex then all the stated inequalities hold but are not strict.  $\blacksquare$

Next, consider the matching  $M$  and the cost-reducing path  $P = (v_1, \dots, v_k)$ . Let  $x = \deg_M(v_1) - 1$  and  $y = \deg_M(v_k)$ . The change in  $\text{cost}_f$  at  $v_1$  due to switching edges along path  $P$  is  $f(x) - f(x+1)$ . Similarly, the change in  $\text{cost}_f$  at  $v_k$  is  $f(y+1) - f(y)$ , therefore

$$\text{cost}_f(M') = \text{cost}_f(M) - (f(x+1) - f(x)) + (f(y+1) - f(y))$$

applying Proposition 2.1.2, and the thesis follows.  $\square$

This result is very important because it provides an extension to Theorem 2.1, which can be reformulated as follows.

**Theorem 2.2** (Optimal Semi-Matching). *Let  $M$  be a semi-matching in  $G$ . If  $\text{cost}_f$  is a strict cost function and if  $M$  is optimal with respect to  $\text{cost}_f$ , then no cost-reducing path relative to  $M$  exists. Conversely, if no cost-reducing path relative to  $M$  exists, then  $M$  is optimal with respect to  $\text{cost}_f$ .*

*Proof.* The proof is the same as the one for Theorem 2.1.  $\square$

This theorem implies a strong corollary that allows us to say that if  $M$  has no cost-reducing paths, then  $M$  is optimal with respect to *every* cost function.

**Corollary 2.2.1.** *A semi-matching  $M$  that is optimal with respect to any strict cost function is optimal with respect to every cost function, strict or weak.*

We are now ready to show that an optimal semi-matching minimizes, among many different objectives, the maximum load on any machine, thus satisfying the primary request of **MRS**.

Let  $X = (x_1, \dots, x_n)$  denote the *load vector* of  $M$ , where  $x_i = \deg_M(v_i)$ . In the following theorem, we prove that an optimal semi-matching is also optimal with respect to every  $L_p$ -norm<sup>1</sup> of the vector  $X$ , for any finite  $p$ .

**Theorem 2.3.** *Let  $1 < p < \infty$ . A semi-matching is optimal if and only if it is optimal with respect to the  $L_p$ -norm of its load vector.*

*Proof.* Fix any  $p > 1$  and define  $f_p(x) = x^p$ . Note that  $f_p$  is strictly convex, so  $\text{cost}_{f_p}$  is a strict cost function. Let  $X$  be the load vector for  $M$ . Since  $\|X\|_p = \text{cost}_{f_p}(M)^{1/p}$ ,  $M$  optimizes  $\|X\|_p$  if and only if it optimizes  $\text{cost}_{f_p}$ . By Corollary 2.2.1,  $M$  is optimal with respect to  $\text{cost}_{f_p}$  if and only if  $M$  is an optimal semi-matching.  $\square$

Our primary objective for **MRS**, however, is to minimize the maximum load on every machine. This is achieved by minimizing the  $L_\infty$ -norm<sup>2</sup> of the load vector  $X$  as we show below.

**Proposition 2.3.1.** *Let  $f_\infty(x) = (n+1)^x$ . A semi-matching that minimizes  $\text{cost}_{f_\infty}$  minimizes the maximal load on a single  $V$ -vertex.*

*Proof.* Let  $M$  be a semi-matching that minimizes  $\text{cost}_{f_\infty}$ . Assume, for the sake of contradiction, that  $M$  does not minimize the maximal load on a single  $V$ -vertex, and that this objective is achieved by a different semi-matching  $M'$ . Let  $\Delta$  and  $\Delta'$  be the maximal load on a single  $V$ -vertex in  $M$  and  $M'$  respectively. Thus,  $\text{cost}_{f_\infty}(M) \geq (n+1)^\Delta$  and  $\text{cost}_{f_\infty}(M') \leq n(n+1)^{\Delta'}$  (this is because, at worst, all the  $n$   $V$ -vertices have load  $\Delta'$ ). However, since  $\Delta' \leq \Delta - 1$ , it follows that  $\text{cost}_{f_\infty}(M') < \text{cost}_{f_\infty}(M)$ , which is a contradiction.  $\square$

This claim leads us to conclude the following.

**Theorem 2.4.** *An optimal semi-matching  $M$  is also optimal with respect to the  $L_\infty$ -norm.*

*Proof.* Since  $\text{cost}_{f_\infty}$  is a cost function, Corollary 2.2.1 implies that  $M$  is optimal with respect to  $\text{cost}_{f_\infty}$ . By Proposition 2.3.1,  $M$  minimizes the maximal load on a  $V$ -vertex and therefore  $M$  is optimal with respect to the  $L_\infty$ -norm.  $\square$

<sup>1</sup>For any  $p \in \mathbb{N}, p > 1$ , the  $L_p$ -norm of vector  $\mathbf{x} = (x_1, \dots, x_k)$  is defined as  $\|\mathbf{x}\|_p = \left( \sum_{i=1}^k |x_i|^p \right)^{1/p}$ .

<sup>2</sup>The  $L_\infty$ -norm of vector  $\mathbf{x}$  is defined as the maximum absolute value of its components.

**Remark.** The converse of Theorem 2.4 is not valid, i.e., minimizing the  $L_\infty$ -norm does not imply minimization of other cost functions. This is because the converse of Proposition 2.3.1 does not hold: minimizing the makespan does not necessarily minimize  $\text{cost}_{f_\infty}$ .

We now have all the theoretical background needed to describe the following two optimal algorithms.

## 2.2 $\mathcal{A}_{SM1}$ : an $\mathcal{O}(mk)$ optimal algorithm

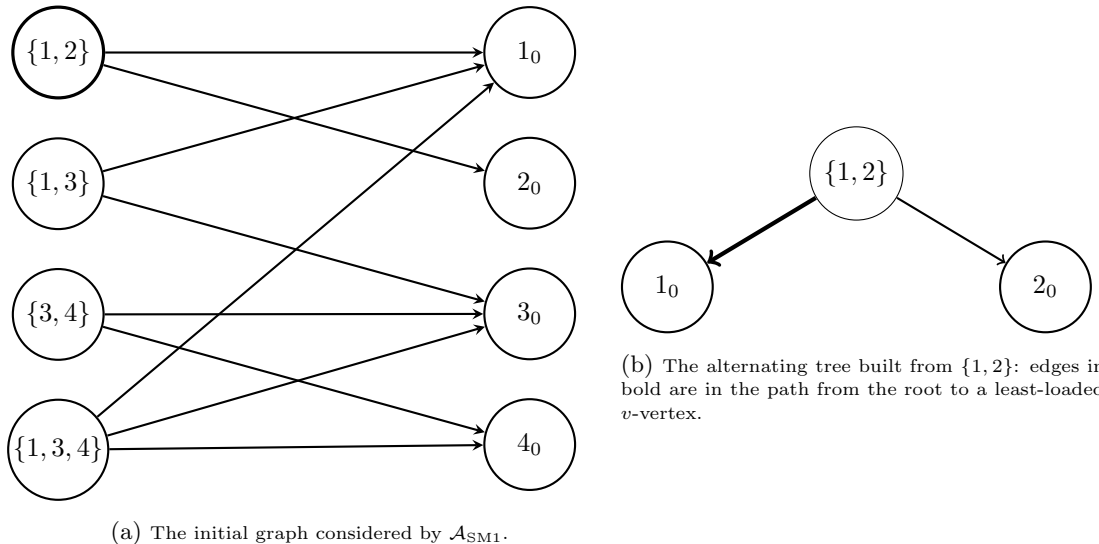
The first algorithm we describe is a modification of the Hungarian algorithm [13] for bipartite matching and it has a time complexity of  $\mathcal{O}(|U||E|)$  or, equally,  $\mathcal{O}(mk)$  in terms of **MRS**.

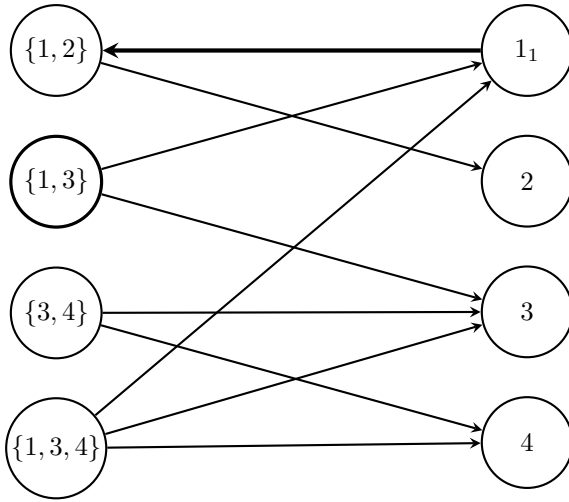
Basically, the Hungarian algorithm for finding maximum bipartite matchings, considers each left-hand vertex  $u$  in turn and builds an *alternating search tree*, i.e. a tree of alternating paths, rooted at  $u$ . When an unmatched right-hand vertex  $v$  (that is a vertex  $v \in V$  with  $\text{deg}_M(v) = 0$ ) is found, the matching and non-matching edges along the  $u$ - $v$  path are switched so that  $u$  and  $v$  are no longer unmatched.

Similarly, this algorithm [4] maintains a partial semi-matching  $M$  starting with the empty set. In each iteration, it considers a vertex  $u \in U$  and builds an alternating search tree from it using a directed breadth-first search. To perform this directed search we give an orientation to the edges in  $G$  according to  $M$ : the edges in  $M$  are directed from  $V$  to  $U$  and edges not in  $M$  are directed from  $U$  to  $V$ . We then select, in this tree, a path  $P$  from  $u$  to the  $v$ -vertex with the least load possible and we form  $P \oplus M$  by switching the orientation of the edges in  $P$ . This will result in  $u$  being no longer unmatched and  $\text{deg}_M(v)$  increased by 1. We repeat this process for every  $u \in U$  and since each iteration matches  $u$  with a single  $v$ -vertex and does not change  $\text{deg}_M(u)$  for any other  $u \in U$ , the resulting selection of edges is a semi-matching.

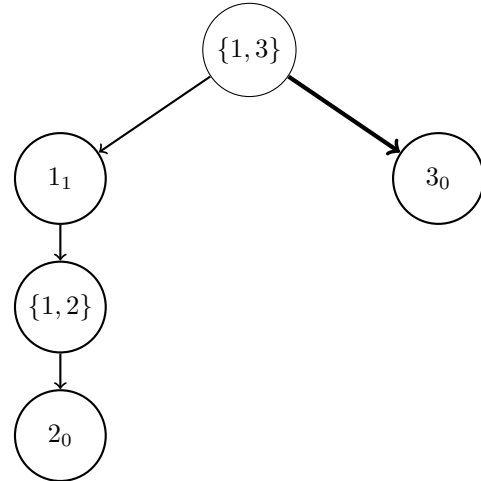
Before discussing the correctness and the time complexity of  $\mathcal{A}_{SM1}$  we provide an example to visualize how the algorithm works. Recall that the input graph of this algorithm is directed.

Figure 2.2: Simulation of  $\mathcal{A}_{SM1}$  on the instance  $\langle \mathcal{N} = \{1, 2, 3, 4\}, \mathcal{F} = \{\{1, 2\}, \{1, 3\}, \{3, 4\}, \{1, 3, 4\}\}$ . The numbers in subscript denote the *load* of the corresponding vertex.

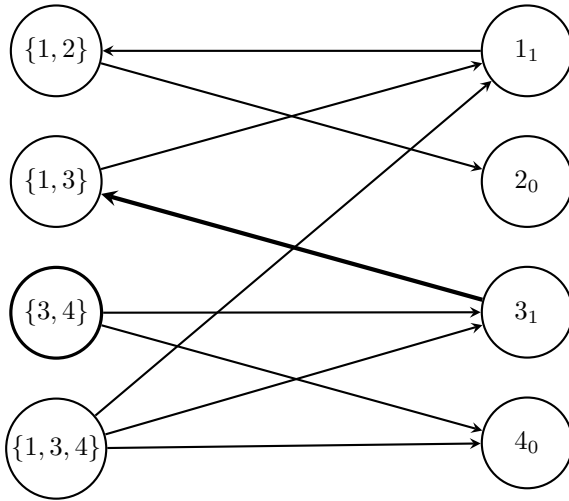




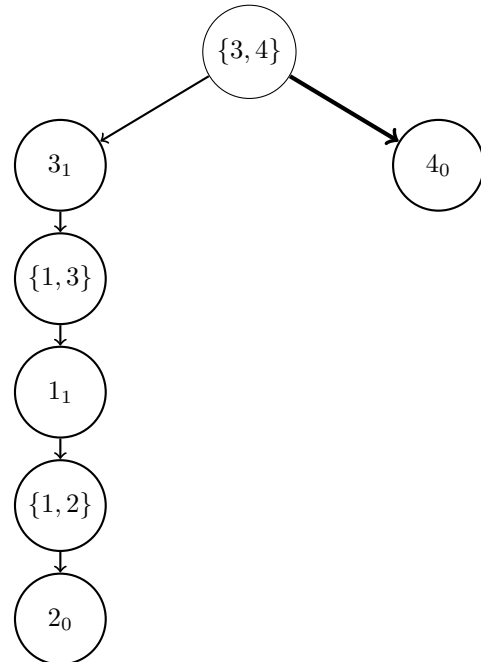
(c) The updated graph obtained by switching the orientation of the bold edges.



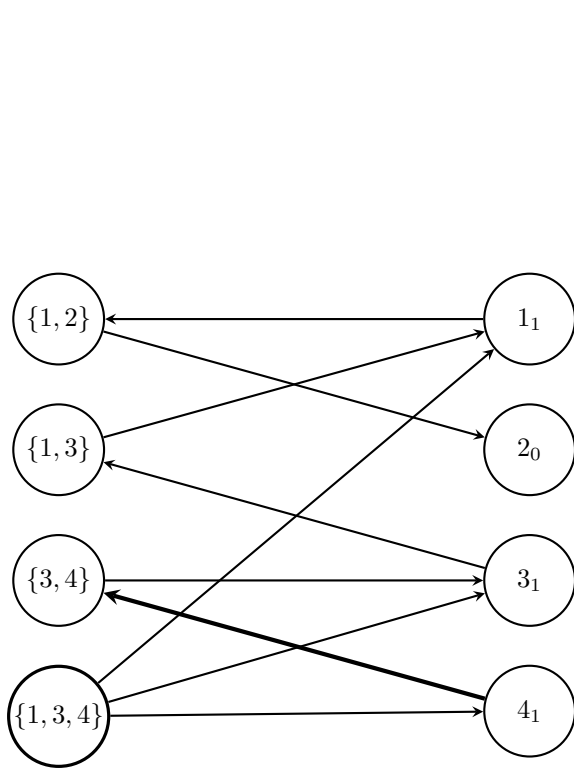
(d) The alternating tree built from  $\{1, 3\}$ .



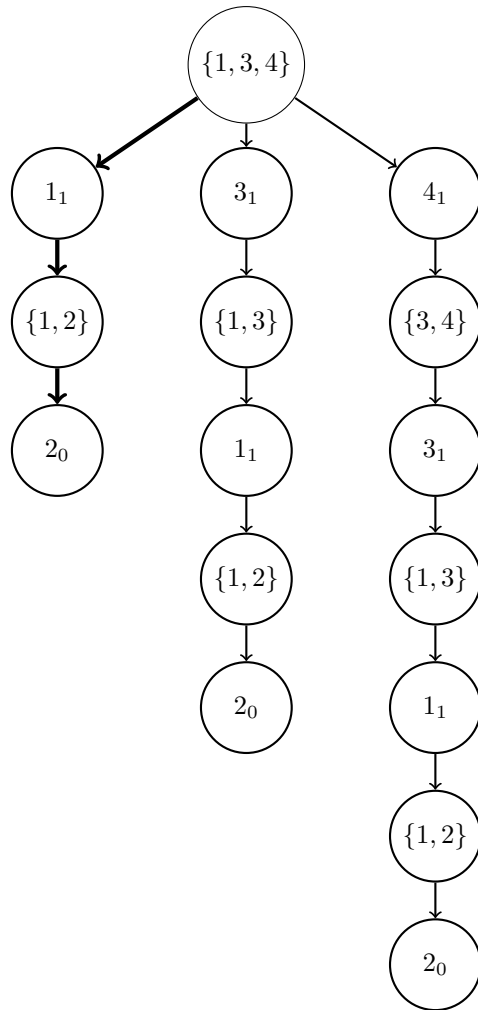
(e) The updated graph.



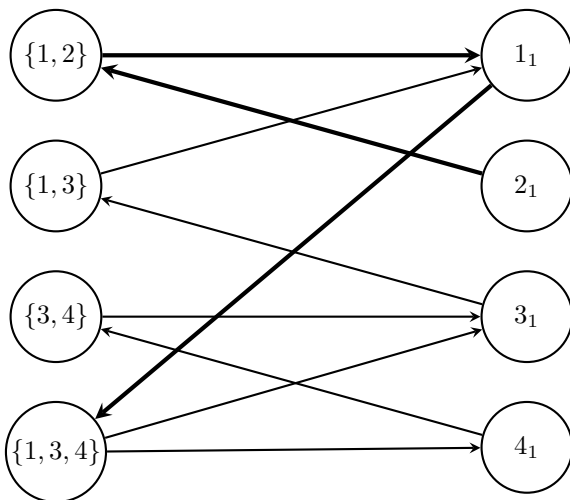
(f) The alternating tree built from  $\{3, 4\}$ .



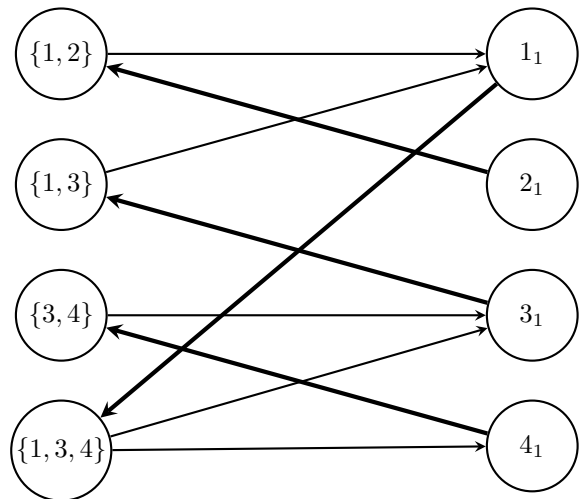
(g) The updated graph.



(h) The alternating tree built from  $\{1, 3, 4\}$ : note that now the least loaded  $v$ -vertex is 2.



(i) The updated graph.



(j) The result computed by  $\mathcal{A}_{SM1}$ . Edges in  $M$  are the ones in bold, directed from right to left. The corresponding solution of **MRS** is then **2, 3, 4, 1**.

### 2.2.1 Correctness and complexity analysis

We now prove the correctness of  $\mathcal{A}_{SM1}$  [4].

**Theorem 2.5.** *Algorithm  $\mathcal{A}_{SM1}$  produces an optimal semi-matching.*

*Proof.* We show that no cost-reducing path is created during the execution of the algorithm. In particular, no cost-reducing path exists at the end of its execution. Thus, by Theorem 2.1, the resulting matching is optimal.

Assume the opposite, and let  $P^* = (v_1, u_1, \dots, v_{k-1}, u_{k-1}, v_k)$  be the *first* cost-reducing path created by the algorithm. Let  $M_1$  be the matching after the iteration in which  $P^*$  is created. Thus,  $\deg_{M_1}(v_1) > \deg_{M_1}(v_k) + 1$ . Without loss of generality, by taking a sub-path of  $P^*$ , we can assume that there exists some  $x$  such that  $\deg_{M_1}(v_1) \geq x + 1$ ,  $\deg_{M_1}(v_i) = x$ ,  $\forall i \in \{2, \dots, k-1\}$  and  $\deg_{M_1}(v_k) \leq x - 1$ . Consider now the last iteration in which the load on  $v_1$  is increased and let  $u'$  be the  $U$ -vertex added to the assignment at this iteration: by definition of the algorithm,  $v_1$  is a least loaded  $V$ -vertex reachable from  $u'$ . Thus, the search tree built for  $u'$  includes only  $V$ -vertices with load at least  $x$ , in particular,  $v_k$  is not reachable from  $u'$ .

Given that the path  $P^*$  exists, at some iteration occurring after the one in which  $u'$  is added, all the edges  $(u_i, v_i)$  of  $P^*$  are in the matching. Let  $u^*$  be the  $U$ -vertex, added after  $u'$ , whose addition to the assignment creates  $P^*$ . We show a contradiction to the way  $u^*$  is assigned. Specifically, we show that when adding  $u^*$ , the algorithm increases the load on some vertex with load at least  $x$ , while  $v_k$ , whose load at that time is at most  $x - 1$  is also reachable from  $u^*$ . To this aim, we first state and prove the following two propositions.

**Proposition 2.5.1.** *When  $u^*$  is added, the load on  $v_k$  is at most  $x - 1$  and  $v_k$  is in the alternating tree rooted at  $u^*$ .*

*Proof.* The load on any  $v$ -vertex can only increase during the algorithm. Since the load on  $v_k$  in  $M_1$  is  $x - 1$  and  $u^*$  is added before  $P^*$  exists,  $x - 1$  is an upper bound on the load on  $v_k$  at the time  $u^*$  is added. To see that  $v_k$  is in  $T^*$  (the search tree built during the assignment of  $u^*$ ) recall the assumption that the path  $P^*$  is created by the assignment of  $u^*$ . Let  $(u_i, v_i)$  be the last edge (i.e. the farthest from  $v_1$ ) of path  $P^*$  that is added to the matching in this iteration. Thus,  $v_i$  is reachable from  $u^*$  and since the edge  $(u_j, v_j)$  is already in the matching  $\forall j, i < j \leq k$ , the suffix of  $P^*$  from  $v_i$  to  $v_k$  must also be in  $T^*$ . ■

**Proposition 2.5.2.** *When  $u^*$  is added, the load on some vertex with load at least  $x$  is increased.*

*Proof.* Suppose the opposite. We show that a cost-reducing path exists *before*  $u^*$  is added, contradicting our choice of  $P^*$  as the first cost-reducing path created by the algorithm. Once again, we use the fact that  $P^*$  is created when  $u^*$  is added. Let  $(u_i, v_i)$  be the first edge (i.e. the closest to  $v_1$ ) of  $P^*$  which is added to the matching at this iteration. Thus, before adding  $u^*$ , the vertex  $v_i$  is reachable from  $v_1$ . Let  $P' = (u^*, \dots, v^*)$  be the path from  $u^*$  to the least loaded vertex in  $T^*$ : note that  $v_i$  must appear in  $P'$ . therefore the path from  $v_1$  to  $v_i$  can be extended to reach  $v^*$  using the suffix of  $P'$  from  $v_i$  to  $v^*$ . All the edges of this suffix are available *before*  $u^*$  is added, since they were all available to  $T^*$ . The load on  $v^*$  must be at least  $x$ , otherwise the above path from  $v_1$  to  $v^*$  is a cost-reducing path which exists before  $P^*$ . ■

Combining these two claims contradicts the execution of the algorithm and therefore  $P^*$  cannot exist. □

To bound the running time of  $\mathcal{A}_{SM1}$  we first observe that there will be exactly  $|U|$  iterations, or equally  $m$  iterations in terms of **MRS**. Each iteration requires at most  $\mathcal{O}(|E|) = \mathcal{O}(k)$  time to build the alternating search tree and at most  $\mathcal{O}(\min\{|U|, |V|\}) = \mathcal{O}(|E|) = \mathcal{O}(k)$  time to switch edges along the alternating path. Therefore the total running time of the algorithm is  $\mathcal{O}(|U||E|) = \mathcal{O}(mk) = \mathcal{O}(m^2n)$ .

## 2.3 $\mathcal{A}_{SM2}$ : a faster-in-practice algorithm

The time bound of the previous algorithm is the same as the Hungarian's one, in fact,  $\mathcal{A}_{SM1}$  is just a variant of it. However, in practice,  $\mathcal{A}_{SM1}$  tends to be slower because of the big size of the alternating trees it builds in each iteration. For this reason, we describe another algorithm from [4] which has a theoretically *worst* time bound but that, in practice, performs much better.

Theorem 2.1 proved that a semi-matching is optimal if and only if the graph does not contain a cost-reducing path.  $\mathcal{A}_{SM2}$  uses that result to find an optimal semi-matching as follows:

1. Start with an initial semi-matching  $M$ ;
2. While there exists a cost-reducing path  $P$ , reduce the cost of  $M$  by switching matching and non-matching edges along  $P$ .

Since each iteration reduces the total cost by an integral amount, the cost can only be reduced a finite number of times, therefore this algorithm must always terminate. Moreover, if the initial semi-matching is nearly optimal, this algorithm terminates after few iterations, thus, finding a good initial  $M$  is important and we start by describing this procedure. Note that  $\mathcal{A}_{SM2}$ , however, does not rely on a good initial semi-matching and it can always find an optimal solution.

*Finding an initial semi-matching.* A simple approach would be to assign each left-hand vertex to a right-hand vertex arbitrarily, however, the following *greedy* algorithm have shown to perform well in practice.

First, the  $U$ -vertices are *sorted* by increasing degree. Each  $U$ -vertex is then considered in turn and assigned to a  $V$ -neighbor with least *load*. In case of a tie, a  $V$ -neighbor with least *degree* is chosen. The purpose of this strategy is to allow more constrained vertices (i.e. the ones with fewer neighbors) to “choose” their matching vertices first, leaving the vertices with more “freedom” (i.e. vertices with higher outdegree) available for the successive steps of the algorithm.

The total time required to find this initial matching is  $\mathcal{O}(|E|) = \mathcal{O}(m)$ , since every edge is examined exactly once, and the sorting can be done using a linear-time sorting algorithm such as counting sort or bucket sort.

*Finding cost-reducing paths.* The critical operation of the  $\mathcal{A}_{SM2}$  algorithm is the method to find cost-reducing paths. A simple approach to find an initial  $V$ -vertex of a cost-reducing path can be to grow a depth-first search tree of alternating paths rooted at  $v$ . To determine if  $G$  has any cost-reducing paths at all, it suffices to perform  $n$  depth-first searches, one for each  $v \in V$ . Unfortunately, this approach performs much redundant work. The following claim helps us to avoid this work.

**Proposition 2.5.3.** *Let  $v, w \in V$  be two vertices such that  $\deg_M(v) \geq \deg_M(w)$  and there is an alternating path from  $v$  to  $w$ . If there is no cost-reducing path starting from  $v$ , then there is no cost-reducing path starting from  $w$ .*

*Proof.* We prove the contrapositive. Since there is an alternating path from  $v$  to  $w$ , any cost-reducing path starting from  $w$  can be extended to a cost-reducing path starting from  $v$ .  $\square$

We can now proceed to find cost-reducing paths in  $G$  as follows. First we find a vertex  $v \in V$  such that

- (i)  $v$  has not been visited by any previous depth-first search;
- (ii) its degree (in the matching  $M$ ) is maximum among all such vertices.

To find such a vertex  $v$  quickly, we keep the non-visited  $V$ -vertices ordered by their load in an array of  $|U| + 1$  *buckets*. Next, we build a depth-first search tree of alternating paths rooted at  $v$  in the usual manner. If this depth-first search tree does not contain a cost-reducing path, then Proposition 2.5.3 allows us to conclude that there is no cost-reducing path starting from any of the vertices that were visited. Since this algorithm considers each edge in  $G$  at most once and there are only  $|U| + 1$  buckets, it finds a cost-reducing path or lack thereof in  $\mathcal{O}(|U| + |E|) = \mathcal{O}(|E|) = \mathcal{O}(k)$  time.

After finding a cost-reducing path, we can reduce the cost of the matching by switching matching and non-matching edges along the path. This step clearly takes  $\mathcal{O}(|E|) = \mathcal{O}(k)$  time as well. Recall that this switching process only affects the load on the first and last vertices on the path.

### 2.3.1 Complexity Analysis

As argued earlier, the initial semi matching can be found in  $\mathcal{O}(|E|) = \mathcal{O}(k)$  time. After this initial step, we iteratively find and remove cost-reducing paths. Finding one takes  $\mathcal{O}(|E|) = \mathcal{O}(k)$  time, while switching matching and non-matching edges along it requires  $\mathcal{O}(\min\{|U|, |V|\}) = \mathcal{O}(|E|) = \mathcal{O}(k)$  time. Thus, the runtime of  $\mathcal{A}_{SM2}$  is  $\mathcal{O}(I \cdot |E|) = \mathcal{O}(I \cdot k)$ , where  $I$  is the number of iterations needed to achieve optimality, so it remains to determine  $I$ . Note that a simple bound of  $I = |U|^2$  may be obtained by observing that the worst possible initial matching has total cost at most  $\mathcal{O}(|U|^2)$  and that each iteration reduces the cost by at least 1.

The following theorem derives, instead, an improved bound of  $\mathcal{O}(\min\{|U|^{3/2}, |U||V|\}) = \mathcal{O}(\min\{m^{3/2}, mn\})$ . However, in real instances of **MRS** it is very difficult for this minimum to be  $m^{3/2}$ , since  $m$  is in general much greater than  $n = |V|$  and can even be exponential in  $n$ . This is why we present a restricted version of the theorem, while the complete one can be found in [4].

**Theorem 2.6.**  $\mathcal{A}_{SM2}$  requires at most  $\mathcal{O}(|U||V|)$  iterations to achieve optimality.

*Proof.* For a given initial semi-matching  $M_0$ , assume that the  $V$ -vertices are sorted by load in non-increasing order (i.e.  $\deg_{M_0}(v_1) \geq \dots \geq \deg_{M_0}(v_m)$ ). In each iteration, the algorithm identifies a cost-reducing path  $P = (v_a, \dots, v_b)$ . By switching matching and non-matching edges along this path we remove one unit of load from  $v_a$  and add one unit to  $v_b$ . This load balancing process can be described as follows: initially we have  $n$  towers of blocks, where the  $i$ th tower consists of  $\deg_{M_0}(v_i)$  blocks. Reducing the cost using path  $P$  amounts to move one block from the tower corresponding to  $v_a$  to the tower corresponding to  $v_b$ . We can bound the number of iterations of  $\mathcal{A}_{SM2}$  by bounding the total number of possible block moves. The proof is based on the following properties:

1. Since we always move load units to a less loaded vertex, operations may be ordered such that the  $V$ -vertices remain sorted in non-increasing order of their loads throughout the execution of the algorithm. In other words, the towers are always sorted from highest to lowest.
2. By the definition of cost-reducing paths, when moving a load unit from  $v_a$  to  $v_b$ , the load on  $v_a$  before the iteration is *strictly* greater than the load on  $v_b$  after the iteration.

For any value of  $n$ , the total number of moves is at most  $nm$  since by the above properties each block may be moved at most  $n$  times.  $\square$

Therefore we can state that the total running time of  $\mathcal{A}_{SM2}$  is  $\mathcal{O}(|U||V||E|) = \mathcal{O}(mnk) = \mathcal{O}(m^2n^2)$  which is worse than the previous algorithm. However, the greedy initialization in practice performs so well that often significantly fewer iterations are required to achieve optimality, making this algorithm the preferred choice. Note that it is possible to create “bad” input graphs, in which the number of iterations needed is  $\Omega(m^{3/2})$ , but most cost-reducing paths in these type of instances are very short, thus each iteration takes roughly constant time.

## Chapter 3

# Heuristics

In this chapter we discuss three heuristics for **MRS**. It may seem quite odd to develop heuristics for a problem that we can solve always *exactly*, however, these three heuristics, and in particular the first two, are much faster than the optimal algorithms described in Chapter 2. For this reason, it may be a good choice to trade an optimal result with a significant reduction of computational time since, as we argued earlier, the size of  $m$  in real instances of **MRS** can be very large. Another reason to develop good heuristics for **MRS**, lies in the  $\mathcal{A}_{SM2}$  initialization step, as it is clearly an initialization with a heuristic. The better this heuristic, the lower iterations of the algorithm needed.

The first two heuristics share many points in common and they are based on a greedy approach. The third heuristic uses instead a completely different, more numeric approach, in a way that will be discussed later in this chapter. From our experiments, we can state that all of the three heuristics have a very high percentage of success, compared with the optimal algorithms.

In the pseudocode of the heuristics we assume to have a field called *representative* in each  $S \in \mathcal{F}$  as the way to associate the representative to its subset, and also a field *cardinality*, already set, that keeps the cardinality of each subset. In this chapter we distinguish between the *occurrences* and the *frequency* of the variables, which are the number of appearances throughout the input instance and the number of times a variable has been chosen as representative, respectively.

### 3.1 Pure Greedy Heuristic

The first and the most basic strategy we present is a simple “pure” greedy algorithm that always chooses the variable that, so far, has been chosen the least amount of times during its process, hence its name. The algorithm breaks ties by selecting the first variable of the tie, but as a general rule we can say that we choose an arbitrary variable. We keep the sets of  $\mathcal{F}$  sorted by their cardinality in ascending order to have more freedom with the choice of variables as the algorithm proceeds.

This heuristic runs in  $\mathcal{O}(k)$  time, therefore it is linear in the size of the input. We provide the pseudocode and the complexity analysis.

---

**Algorithm 2** (Pure Greedy) Always selects the least-chosen variable.

---

```

1: function GREEDY(Sets)
2:   SORT(Sets)                                ▷ according to their cardinality in ascending order
3:   freq  $\leftarrow$  new array of length n
4:   for i  $\leftarrow$  1 to n do
5:     freq[i]  $\leftarrow$  0
6:   for all S  $\in$  Sets do
7:     min  $\leftarrow$   $\infty$ 
8:     repr  $\leftarrow$  -1
9:     for all v  $\in$  S do
10:      if freq[v] < min then
11:        repr  $\leftarrow$  v
12:        min  $\leftarrow$  freq[v]
13:     S.representative  $\leftarrow$  repr
14:     freq[repr]++

```

---

Consider, for example, the same input instance as the one of Section 1.1,

$$\langle \mathcal{N} = \{1, 2, 3, 4, 5\}, \mathcal{F} = \{ \{1, 2\}, \{1, 3\}, \{2, 3\}, \{2, 4\}, \{4, 5\} \} \rangle.$$

The PURE GREEDY algorithm produces the following selection of representatives: 1, 3, 2, 4, 5 (the order is the same of the subsets in the input), which is an optimal solution with maximum frequency of 1. However, with the following instance things get worse.

$$\langle \mathcal{N} = \{1, 2, 3, 4\}, \mathcal{F} = \{ \{1, 2\}, \{1, 3\}, \{3, 4\}, \{1, 3, 4\} \} \rangle$$

The algorithm produces the result 1, 3, 4, 1, with a value of 2. This is clearly not an optimal solution, because the optimal value for this instance is 1 and is achieved, for example, by the choice of 2, 1, 3, 4. Specifically, the heuristic here fails because of the tie-breaking mechanism that chose 1 in the first subset instead of 2, and it is not a problem of the implementation. Suppose that we break ties by choosing the *second* variable in the tie instead of the first: it is clear that this instance will be correctly solved. But if we swap the first and the second subset of the input, we get back with the same problem<sup>1</sup>. This happens because the algorithm *is not aware* of what the subsequent sets will be like, and blindly proceeds throughout them.

**Theorem 3.1.** *The running time of the pure greedy heuristic is  $\mathcal{O}(k) = \mathcal{O}(mn)$ .*

*Proof.* We can sort the sets in  $\mathcal{O}(m + n)$  time using a linear time sorting algorithm. This is possible because we sort the sets according to their cardinality, thus the sorting merely becomes a sorting of integers bounded by the value of  $n$ , which is the highest cardinality possible for each subset in  $\mathcal{F}$ . The second **for** loop is just a complete scan of the input, therefore it takes  $\mathcal{O}(k)$  time. The total running time is then  $\mathcal{O}(m + n + k) = \mathcal{O}(k) = \mathcal{O}(mn)$ .  $\square$

---

<sup>1</sup>Here we are supposing that the linear-time sorting algorithm is also stable. Such an algorithm is, for example, *Counting Sort*.

As a practical consideration, we can say that often  $k \ll mn$ , as the  $mn$  limit is reached only when *all* of the  $m$  subsets are of cardinality  $n$ , which is an extreme situation that almost never appears in real instances of MRS.

## 3.2 Forward Greedy Heuristic

We developed this heuristic with the problem of *blindness* that affects the pure greedy one in mind, thus we decided to make the algorithm follow the opposite direction or, in other words, to focus primarily on the future occurrences of the variables and not only on the past frequencies. For this reason, we say that this heuristic uses a “forward greedy” approach. The sets are kept sorted by their cardinality as in the previous heuristic for the same reason. More formally, the selection strategy of this algorithm is as follows:

1. Among all the variables in the current subset that have not been chosen more than  $\lceil \frac{m}{n} \rceil$  times, choose the one with lowest number of occurrences across all the (remaining) subsets. If such a variable does not exist, choose a variable using the pure greedy strategy;
2. In case of ties, select the variable that has been chosen the lowest number of times, i.e. the variable with the least frequency (as in the pure greedy heuristic);
3. In case of successive ties, select a variable arbitrarily.

As before, we provide the pseudocode of this heuristic and we discuss its time complexity.

---

### Algorithm 3 (Forward Greedy)

---

```

1: function FORWARDGREEDY(Sets)
2:   SORT(Sets)                                ▷ according to their cardinality in ascending order
3:   freq ← new array of length  $n$ 
4:   occ ← new array of length  $n$                 ▷ contains the occurrences of each variable
5:   for  $i \leftarrow 1$  to  $n$  do
6:     freq[ $i$ ] ← 0
7:   for all  $S \in$  Sets do
8:     for all  $v \in S$  do
9:       occ[ $v$ ]++                               ▷ count the occurrences of each variable
10:  for all  $S \in$  Sets do
11:    candidatef ← the first variable in  $S$  with the lowest freq value
12:    candidateo ← the first variable in  $S$  with the lowest value of occ
13:    if freq[candidateo] <  $\lceil \frac{m}{n} \rceil$  then      ▷ if we already chose candidateo less than  $\lceil \frac{m}{n} \rceil$  times
14:      S.representative ← candidateo
15:    else
16:      S.representative ← candidatef
17:
18:    for all  $v \in S$  do
19:      occ[ $v$ ] --                               ▷ remove the occurrences of all the variables contained in  $S$ 
20:
21:    freq[S.representative]++

```

---

The reason of the check at line 13 is based on Proposition 1.2.1: in an ideal solution, its value would be  $\lceil \frac{m}{n} \rceil$  and so we try to not exceed this mean when choosing a variable. Note that we use this lower bound as a basic approach because the possibly-higher lower bound provided in Theorem 1.3 takes exponential time to be calculated.

As an example, consider the second input instance provided in the previous Section:

$$\langle \mathcal{N} = \{1, 2, 3, 4\}, \mathcal{F} = \{ \{1, 2\}, \{1, 3\}, \{3, 4\}, \{1, 3, 4\} \} \rangle.$$

This time, the heuristic gives the solution 2, 1, 3, 4 with value 1, which is optimal. Note that the result does not depend on the order of the subsets in the input. However, this is not true in general.

**Theorem 3.2.** *The running time of the forward greedy heuristic is  $\mathcal{O}(k) = \mathcal{O}(mn)$ .*

*Proof.* We can sort the sets using the same linear-time sorting algorithm as in the pure greedy heuristic, at the cost of  $\mathcal{O}(m + n)$  time. The first **for all** loop is a scan of the input and it takes  $\mathcal{O}(k)$  time. Finding the variable with lowest frequency and the variable with lowest occurrences can be done with a scan of the current  $S$ , which takes  $\mathcal{O}(n)$  time, and the same argument applies to the last **for all** loop. Therefore, the **for all** loop starting at line 10 takes  $\mathcal{O}(k)$  time. Therefore the total running time is  $\mathcal{O}(m + n + k) = \mathcal{O}(k) = \mathcal{O}(mn)$ .  $\square$

We finally note that the initialization step of  $\mathcal{A}_{\text{SM2}}$  is a hybrid between our pure greedy heuristic and this one, because it chooses according only to *frequencies* (referred as the *load* in  $\mathcal{A}_{\text{SM2}}$ ), but breaks ties by choosing the variable with the least *occurrences* (previously referred as the *degree*).

### 3.3 Pagerank inspired Heuristic

The last heuristic we describe uses a completely different strategy to build a solution, and it is a more numerical one compared with the two previous heuristics.

Unlike any other previous algorithm (including optimal ones), this algorithm creates a solution by iteratively *removing* variables from the subsets, i.e., at each iteration it excludes a variable from a subset. At the end of its execution there will be only one variable left in each subset and they will be chosen as representatives. For this reason, we chose the name “Pagerank” for the heuristic, according to the nature of the Markov Matrix used by Google, that is obtained by successive transformations of the initial hyperlinking matrix [14].

We build a  $m \times n$  matrix  $W$  (called the *relative weight matrix* or *weight matrix*) with variables on the columns and subsets on the rows and we populate it with the uniform probability that we have, at any given time, to choose the variable  $j$  from the  $i$ -th subset in  $\mathcal{F}$ .

**Definition 3.1** (Weight Matrix). The *weight matrix* is a  $m \times n$  matrix,  $W$ , with

$$W_{i,j} = \begin{cases} \frac{1}{|\mathcal{F}[i]|} & \text{if } j \in \mathcal{F}[i] \\ 0 & \text{otherwise} \end{cases}.$$

We define the *total weight* (sometimes only *weight* in the following) of a variable as the sum over its column, that is,  $weight(j) = \sum_{i=1}^m W_{i,j}$ . Our goal is to minimize the weight of each variable as much as possible, in particular the maximum one. The algorithm proceeds as follows:

1. Choose the variable with maximum weight, and look for a subset that contains it with its highest relative weight (without considering the “1”s). In case of ties on the maximum weight, select one arbitrarily, ignoring the variables that have only “1”s in their column. In case of ties among the subsets, select one arbitrarily;
2. Exclude the variable from the chosen subset by setting its relative weight to 0;
3. Update the relative weights of the variables in the chosen subset, having its cardinality *logically* decreased by one;
4. Repeat this procedure until there is only one variable left in each subset, i.e., in every row of the matrix there are  $n - 1$  zeros and one “1”.

Tables 3.1 and 3.2 show an example of iteration of this algorithm.

Subset	1	2	3	4
{1, 2}	0.5	0.5	—	—
{1, 3}	0.5	—	0.5	—
{3, 4}	—	—	0.5	0.5
{1, 3, 4}	0.33	—	0.33	0.33
<b>Total Weights</b>	1.33	0.5	1.33	0.83

Table 3.1: The initial matrix built by the algorithm for the input instance used in the previous sections, ( $\mathcal{N} = \{1, 2, 3, 4\}$ ,  $\mathcal{F} = \{\{1, 2\}, \{1, 3\}, \{3, 4\}, \{1, 3, 4\}\}$ ). The circled numbers denote the maximum relative and total weight (with ties broken arbitrarily).

Subset	1	2	3	4
{1, 2}	—	1	—	—
{1, 3}	0.5	—	0.5	—
{3, 4}	—	—	0.5	0.5
{1, 3, 4}	0.33	—	0.33	0.33
<b>Total Weights</b>	0.83	1	1.33	0.83

Table 3.2: The result of zeroing the variable “1” from the first subset. The new maximum relative and total weights are circled.

Intuitively, prior to the iteration on Table 3.1, we had a probability of  $\frac{1}{2}$  of picking either 1 or 2 from the first subset. After that iteration, the variable 1 *certainly* will not be chosen as representative, with 2 being picked instead (because it has a probability of 1).

We now provide the pseudocode for this heuristic and we discuss its time complexity.

---

**Algorithm 4** (Pagerank)

---

```

1: function PAGERANK(Sets)
2:    $W \leftarrow$  new  $m \times n$  matrix
3:   total_weights  $\leftarrow$  new array of length  $n$ 
4:   for all  $S \in$  Sets do                                      $\triangleright$  Initialization of the matrix
5:     for all  $v \in S$  do
6:        $W[S][v] \leftarrow \frac{1}{|S|}$ 
7:       total_weights[ $v$ ]  $\leftarrow$  total_weights[ $v$ ] +  $W[S][v]$ 
8:   for  $i \leftarrow 1$  to  $k - m$  do
9:      $v \leftarrow$  the variable with maximum total_weight
10:     $S \leftarrow$  a subset in which  $v$  appears with highest relative weight
11:    total_weights[ $v$ ]  $\leftarrow$  total_weights[ $v$ ] -  $W[S][v]$ 
12:     $W[S][v] \leftarrow 0$ 
13:    Decrease the (logical) cardinality of  $S$  by 1
14:    for all  $x \in S \setminus \{v\}$  do
15:       $W[S][x] \leftarrow W[S][x] + \frac{1}{|S|(|S|+1)}$             $\triangleright$  With  $|S|$  being the logical cardinality of  $S$ 
16:      total_weights[ $v$ ]  $\leftarrow$  total_weights[ $v$ ] +  $\frac{1}{|S|(|S|+1)}$ 
17:   for all  $S \in$  Sets do
18:     S.representative  $\leftarrow$  the only variable of  $S$  with  $W[S][variable] = 1$ 

```

---

**Theorem 3.3.** *The running time of the pagerank heuristic is  $\mathcal{O}(mn^2 \log n)$ .*

*Proof.* We use  $n$  balanced binary search trees (we use AVL trees in our implementation) to keep track of which variable is in which subset with a particular relative weight. Specifically, for each variable in  $\mathcal{N}$  we keep a pointer to the root of its AVL, with its nodes consisting of:

- A key, which is one of the possible  $n - 2$  relative weights (we do not keep track of the weights 0 and 1);
- A hash table,  $H$ , of length  $m$ ;
- An array,  $A$ , of length  $m$ ;
- The current size of the array  $A$  (and therefore the index of the first empty slot in  $A$ ), namely  $A.size$ .

The satellite data keep track of which subsets contain the variable represented by the AVL. Note that a list, a hash table or an array alone do not allow us to make insertion, deletions, and search of elements in constant time, a result that can be instead achieved with a pair of a hash table and an array using the following techniques:

---

```

1: function INSERT(set)
2:    $A[A.size] \leftarrow set$ 
3:    $H[set] \leftarrow A.size$ 
4:    $A.size++$ 

```

---



---

```

1: function REMOVE(set)
2:    $d \leftarrow A[A.size - 1]$ 
3:    $i \leftarrow H[set]$ 
4:    $A[i] \leftarrow d$ 
5:    $H[d] \leftarrow i$ 
6:    $H[set] \leftarrow \text{NIL}$ 
7:    $A.size --$ 

```

---

The number of iterations of the algorithm is fixed to be *exactly*  $k - m$ , since we zero a variable from a subset at each iteration, leaving one variable in every subset, thus zeroing exactly  $k - m$  variables in total.

At each iteration, we can find the maximum total weight in constant time by keeping it updated every time we modify a relative weight. Finding a subset in which the variable to be zeroed appears with its highest relative weight can be done in  $\mathcal{O}(n \log n)$  time in the following way:

- (i) In the corresponding AVL, get the node with the maximum key (requires  $\mathcal{O}(\log n)$  time);
- (ii) If that node contains at least one element, i.e., if  $A.size > 0$ , grab an element (for example the first, or also a random one) from the node (requires  $\mathcal{O}(1)$  time);
- (iii) If that node does not contain any elements, repeat the process with the predecessor node (requires at most  $n$  searches of a predecessor, each of them takes  $\mathcal{O}(\log n)$  time  $\implies \mathcal{O}(n \log n)$  total time).

Zeroing the variable from the subset found above takes  $\mathcal{O}(1)$  time, and updating the relative weight of the other variables in that subset (and their total weights) requires  $\mathcal{O}(n)$  time. To update the corresponding AVLs we may simply remove and re-insert the subsets' index in them, requiring  $\mathcal{O}(\log n)$  total time. Therefore, each iteration takes  $\mathcal{O}(n \log n)$  time, and since there will be  $k - m$ , or equally  $\mathcal{O}(k)$ , iterations we can conclude that the algorithm runs in  $\mathcal{O}(kn \log n) = \mathcal{O}(mn \cdot n \log n) = \mathcal{O}(mn^2 \log n)$  time.  $\square$

The time bound just obtained is higher than the ones for the previous heuristics. However, thanks to the nature of **MRS**, we know that almost always  $n \ll m$ , therefore this time bound will be practically much lower than the one of  $\mathcal{A}_{\text{SM1}}$ , which is  $\mathcal{O}(m^2 n)$ .

A refinement for this heuristic could be to not break ties arbitrarily when selecting the subset with the variable of highest relative weight, but rather to choose the subset that contains that variable *and* the variable of *lowest* total weight. This is because the weights of the variables will denote, at the end of

the algorithm, the number of times they have been chosen. Zeroing a variable with high total weight reflects the attempt to minimize the maximum frequency, and doing so from a subset that contains also a variable much “lighter” will increase its total weight, which is the least dangerous at the moment. If such a subset does not exist, look for the *second* lowest total weight and so on. However, we were not able to stay under the  $\mathcal{O}(m^2n)$  time limit with this refinement, and in this case the algorithms described in Chapter 2 are a superior choice.



## Chapter 4

# Experimental Results

In this chapter we provide some statistics about the practical execution time of all the algorithms presented in the previous chapters along with statistics on the quality of the three heuristics, obtained by comparing their result to the one of the optimal algorithms. To represent the instances of **MRS** as close to the reality as possible (i.e., to use numbers that can actually be involved during the process of synthesis on switching lattices of Boolean functions), we tested the algorithms under the assumption that  $m \gg n$ . For example, for every value of  $n$ , we will generate instances with,  $m = \Theta(n^2)$ ,  $\Theta(n^3)$  and, when possible,  $m = 2^{n-1}$ .

The test cases can be divided into two large groups, based on the probability distribution of the variables across the subsets: we chose a Normal distribution and an Uniform distribution. Roughly speaking, one can expect that the instances generated with the uniform probability will be harder than the ones generated with the Gaussian one, because each variable will tend to appear the same amount of times across the subsets, making any kind of simplification (like the one described in Corollary 1.3.1) nearly impossible. The cardinality of the subsets always follows the uniform distribution. Table 4.1 and 4.2 are organized as follows:

1. The first column contains the pair  $\langle n, m \rangle$ . For each of these pairs, we generate 100 independent instances;
2. The second column contains the average  $k$ , which corresponds to the average size of the input across the 100 instances. More formally,  $k_{avg} = \frac{\sum_{i=1}^{100} \sum_{j=1}^m |\mathcal{F}_i[j]|}{100}$ ;
3. The remaining columns represent the average time spent by each algorithm to compute the solution. Like before, it is obtained by averaging the time spent for each of the 100 input instances;
4. The second table shows the same kind of data for input instances generated with the uniform distribution.

$\langle n, m \rangle$	$k_{avg}$	$\mathcal{A}_{SM1}$	$\mathcal{A}_{SM2}$	Greedy	GreedyFWD	Pagerank
$\langle 10, 100 \rangle$	515	0.005371s	0.000065s	0.000017s	0.000024s	0.000309s
$\langle 10, 512 \rangle$	2756	0.142710s	0.000292s	0.000070s	0.000113s	0.001240s
$\langle 20, 400 \rangle$	4198	0.174935s	0.000367s	0.000088s	0.000164s	0.004596s
$\langle 10, 1000 \rangle$	5464	0.578575s	0.000484s	0.000124s	0.000194s	0.002187s
$\langle 30, 900 \rangle$	13958	1.947167s	0.000936s	0.000184s	0.000400s	0.031149s
$\langle 50, 2500 \rangle$	63809	24.948925s	0.003540s	0.000683s	0.001695s	0.272875s
$\langle 20, 8000 \rangle$	83976	103.82581s	0.005931s	0.001307s	0.002416s	0.107823s
$\langle 30, 27000 \rangle$	418418	1725.06s	0.024287s	0.004805s	0.009786s	0.830832s
$\langle 50, 125000 \rangle$	3188839	†	0.162482s	0.033316s	0.076941s	15.668205s
$\langle 20, 2^{19} \rangle$	5504172	†	0.371158s	0.088621s	0.168090s	8.491898s

Table 4.1: Average execution times over 100 randomly-generated instances with the Gaussian distribution. Results are marked with a '†' when an algorithm is stopped after 18 hours of computation.

$\langle n, m \rangle$	$k_{\text{avg}}$	$\mathcal{A}_{\text{SM1}}$	$\mathcal{A}_{\text{SM2}}$	Greedy	GreedyFWD	Pagerank
$\langle 10, 100 \rangle$	552	0.005762s	0.000072s	0.000019s	0.000028s	0.000350s
$\langle 10, 512 \rangle$	2819	0.161303s	0.000289s	0.000057s	0.000109s	0.001537s
$\langle 20, 400 \rangle$	4198	0.178702s	0.000313s	0.000065s	0.000148s	0.004927s
$\langle 10, 1000 \rangle$	5491	0.684278s	0.000540s	0.000134s	0.000220s	0.002989s
$\langle 30, 900 \rangle$	13949	2.005043s	0.000825s	0.000173s	0.000431s	0.029007s
$\langle 50, 2500 \rangle$	63739	23.337943s	0.002764s	0.000581s	0.001639s	0.253003
$\langle 20, 8000 \rangle$	83962	87.022895s	0.004335s	0.001061s	0.002184s	0.095918s
$\langle 30, 27000 \rangle$	418429	†	0.020709s	0.004799s	0.011405s	0.908683s
$\langle 50, 125000 \rangle$	3187990	†	0.134957s	0.029164s	0.079760s	14.569035s
$\langle 20, 2^{19} \rangle$	5504567	†	0.375886s	0.096140s	0.194809s	9.670038s

Table 4.2: Average execution times over 100 randomly-generated instances with the Uniform distribution. Results are marked with a '†' when an algorithm is not executed.

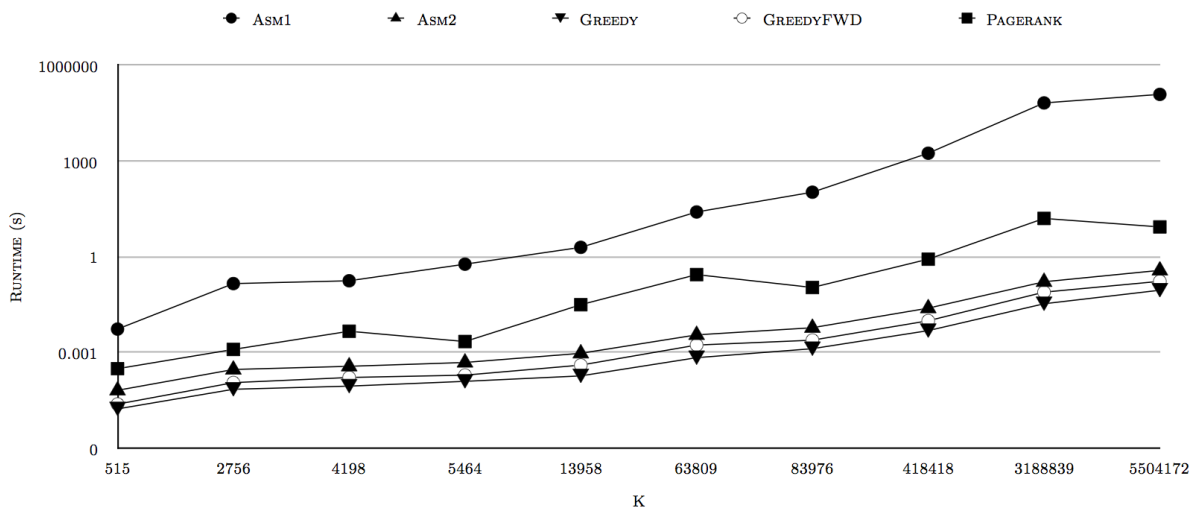


Figure 4.1: Visual comparison of the running times obtained with the Gaussian distribution. The graph obtained with the Uniform distribution is analogous.

These results show how  $\mathcal{A}_{\text{SM1}}$  is practically unfeasible for real instances of **MRS**, since its running time increases drastically even for relatively small increases of the input size, while the other algorithms have no problems in handling even a huge amount of sets in practice, such as  $2^{19}$ . We can note that the average  $k$ , i.e., the average size of each input instance generated, is roughly equal to  $mn/2$ . This is important for the linear-time heuristics, since their running time is practically halved with respect to the  $\mathcal{O}(mn)$  time bound. Another thing to note is that for a major part of Table 4.2 the running times are higher than the ones of Table 4.1, practically confirming our hypothesis that the instances generated with a uniform distribution are generally harder.

We also tested the success rate of our heuristics, comparing their result with the one obtained by  $\mathcal{A}_{\text{SM2}}$ . In the table below we provide the percentage of *failures* of the heuristics, along with their average deviation from the optimal value. We wanted also to test the quality of the  $\mathcal{A}_{\text{SM2}}$ 's initialization heuristic, to compare it against ours. The test cases have the same  $m$  and  $n$  of the tables before, but we also added a few more, smaller, ones.

$\langle \mathbf{n}, \mathbf{m} \rangle$	$k_{\text{avg}}$	Greedy	GreedyFWD	Pagerank	$\mathcal{A}_{\text{SM2}}$ heuristic
$\langle 15, 20 \rangle$	157	—	—	—	—
$\langle 10, 30 \rangle$	162	—	—	—	—
$\langle 10, 50 \rangle$	266	1 (+1)	—	—	—
$\langle 20, 50 \rangle$	505	—	—	—	—
$\langle 10, 100 \rangle$	545	—	—	—	—
$\langle 50, 100 \rangle$	2543	10 (+1)	—	4 (+1)	—
$\langle 10, 512 \rangle$	2752	—	—	—	—
$\langle 20, 400 \rangle$	4203	—	—	—	—
$\langle 100, 100 \rangle$	5060	19 (+1)	—	4(+1)	—
$\langle 10, 1000 \rangle$	5459	—	—	—	—
$\langle 30, 900 \rangle$	13927	—	—	—	—
$\langle 50, 2500 \rangle$	63876	—	—	—	—
$\langle 20, 8000 \rangle$	83930	—	—	—	—
$\langle 30, 27000 \rangle$	418305	—	—	—	—
$\langle 50, 125000 \rangle$	3188873	—	—	—	—
$\langle 20, 2^{19} \rangle$	5504229	—	—	—	—

Table 4.3: Average errors and average deviation from the optimal value of the heuristics over 100 randomly-generated instances with the Gaussian distribution.

$\langle \mathbf{n}, \mathbf{m} \rangle$	$k_{\text{avg}}$	Greedy	GreedyFWD	Pagerank	$\mathcal{A}_{\text{SM2}}$ heuristic
$\langle 15, 20 \rangle$	158	1 (+1)	—	—	—
$\langle 10, 30 \rangle$	158	—	—	—	—
$\langle 10, 50 \rangle$	262	1 (+1)	—	1 (+1)	—
$\langle 20, 50 \rangle$	515	—	—	—	—
$\langle 10, 100 \rangle$	516	—	—	—	—
$\langle 50, 100 \rangle$	2525	7 (+1)	—	5 (+1)	—
$\langle 10, 512 \rangle$	2785	—	—	—	—
$\langle 20, 400 \rangle$	4213	—	—	—	—
$\langle 100, 100 \rangle$	4975	3 (+1)	—	—	—
$\langle 10, 1000 \rangle$	5467	—	—	—	—
$\langle 30, 900 \rangle$	13994	—	—	—	—
$\langle 50, 2500 \rangle$	63894	—	—	—	—
$\langle 20, 8000 \rangle$	83993	—	—	—	—
$\langle 30, 27000 \rangle$	418457	—	—	—	—
$\langle 50, 125000 \rangle$	3188750	—	—	—	—
$\langle 20, 2^{19} \rangle$	5504227	—	—	—	—

Table 4.4: Average errors and average deviation from the optimal value of the heuristics over 100 randomly-generated instances with the Uniform distribution.

As we can see from the tables above, every heuristic performs very well in general with only few errors for relatively small input instances (i.e. for relatively small values of  $k$ ). As the size of the input grows, errors will practically tend to zero. In particular, our Greedy Forward heuristic and the heuristic used for the initialization of  $\mathcal{A}_{\text{SM2}}$  show a 100% success rate which is also extremely useful for the running time of the latter algorithm. In fact, the running time of  $\mathcal{A}_{\text{SM2}}$  is, in practice, twice as the running time of the Greedy Forward heuristic which is in turn very similar to the  $\mathcal{A}_{\text{SM2}}$ 's. This is because checking if the solution obtained by the heuristic does not contain a cost-reducing path is done in  $\mathcal{O}(|E|) = \mathcal{O}(k)$  time (which corresponds to an iteration of the algorithm) and the two heuristics run in  $\mathcal{O}(k)$  time as well. However, we were able to find some input instances that make the latter heuristic failing to catch

the optimal solution. One of these counterexamples is the following:

$$\langle \mathcal{N} = \{1, 2, 3, 4\}, \mathcal{F} = \{ \{1, 2\}, \{1, 3\}, \{1, 4\}, \{3, 4\} \} \rangle$$

The  $\mathcal{A}_{\text{SM2}}$ 's initialization heuristic provides a solution of 2, 3, 4, 3 (having a value of 2). This happens because the heuristic does not update the degree of the  $v$ -vertex it creates during its execution, i.e., does not decrease the number of *future occurrences* of the variables once a subset has been processed. This is in contrast to our Greedy Forward heuristic which, thanks to this updates of the future occurrences, correctly gives the solution 2, 1, 4, 3 (with value 1). However, the tables above show that the  $\mathcal{A}_{\text{SM2}}$ 's initialization heuristic runs in this kind of errors only for very small values of  $n$  and  $m$ , in practice.

## Chapter 5

# Conclusions and Future Work

In the last years, switching lattices have seen renewed interest especially thanks to the increased use of nanoscale technologies in building new electronic *chips* and our work addresses one of the problems that one encounters when developing one such device.

We have seen how our optimization problem can be solved by reducing it to, at least, three different types of problems, each of them known to be polynomially solvable. We chose to focus primarily on the scheduling area, but the other two areas can be explored as well. In fact, it is also possible to generalize the graph orientation problem to hypergraphs<sup>1</sup>, which represent an immediate translation of the instances of **MRS**, just by putting  $V = \mathcal{N}$  and  $E = \mathcal{F}$ . **MRS** can then be solved by directing the hyperedges of the hypergraph (see, for example, [15] to have an idea on how a hypergraph can be directed). The time bound for solving this problem using the network flow reduction, which we recall it to be  $\mathcal{O}(m^3 \log m)$ , is loose. Specifically, the  $m^3$  term comes from the complexity of the algorithm used for solving the maximum flow problem and therefore, the total running time can be reduced by using a more efficient algorithm for solving that problem.

The first of the two optimal algorithm that we described works in the “good old fashion” of constructing a solution starting from an empty one, while the second one improves a starting, possibly already optimal, solution. This shows how “informed” algorithm can beat a theoretically-faster algorithm in practice, and also gives a reason to develop good heuristics. We believe that the Pagerank heuristic, in particular, can serve as a basis for creating a new optimal algorithm that works using a totally different and more numerical approach. For example, one can provide the algorithm with some form of autocorrection, like the one implicitly performed by  $\mathcal{A}_{SM1}$  at the moment of following a path on the alternating search tree.

All of the heuristics we presented can also find application in areas where the approach *has to be* heuristical, thanks to the reductions we have shown in Chapter 1. For example, having a weight on our subsets is equivalent to have a different execution time for each job in a scheduling environment. Moreover, weights on the subsets become weights on the edges of a graph that has to be directed, like in the MMO problem, which we recall it to be the problem of finding an orientation for all edges of an undirected graph such that the maximum outdegree (or equivalently, indegree) is minimum. These problems, when weighted with weights greater than 1, become  $\mathcal{NP}$ -hard, and require heuristic or approximation algorithms.

Again in the context of switching lattices, it can also be a good idea to assign the same variable to a group of neighboring cells to simplify the resulting architecture (and perhaps meeting some practical needs such as cost reductions), and this can be reflected in **MRS**, for example, by putting weights on the variables instead of on the sets, and one can investigate on how all the algorithms proposed here react to weighted instances.

---

<sup>1</sup>An hypergraph  $G = (V, E)$  is a graph in which the edges are formed by subsets of  $V$  of any cardinality, not just pairs of vertices.



# Bibliography

- [1] S. B. Akers, “A rectangular logic array”, *IEEE Transactions on Computers*, vol. 21, no. 8, pp. 848–857, 1972.
- [2] M. Altun and M. D. Riedel, “Logic synthesis for switching lattices”, *IEEE Transactions on Computers*, vol. 61, no. 11, pp. 1588–1600, 2012.
- [3] G. Gange, H. Søndergaard, and P. J. Stuckey, “Synthesizing optimal switching lattices”, *ACM Transactions on Design Automation of Electronic Systems*, vol. 20, pp. 1–14, nov 2014.
- [4] N. J. Harvey, R. E. Ladner, L. Lovász, and T. Tamir, “Semi-matchings for bipartite graphs and load balancing”, *Journal of Algorithms*, pp. 53–78, 2006.
- [5] Y. Asahiro, E. Miyano, H. Ono, and K. Zenmyo, “Graph orientation algorithms to minimize the maximum outdegree”, *International Journal of Foundations of Computer Science*, pp. 197–215, 2010.
- [6] Y. Lin and W. Li, “Parallel machine scheduling of machine-dependent jobs with unit-length”, *European Journal of Operational Research*, pp. 261–266, 2004.
- [7] V. Venkateswaran, “Minimizing maximum indegree”, *Discrete Applied Mathematics*, vol. 143, pp. 374–378, 2004.
- [8] E. Mokotoff, “Parallel machine scheduling problems: A survey”, *Asia-Pacific Journal of Operational Research*, no. 18, pp. 193–242, 2001.
- [9] T. C. E. Cheng and C. C. S. Sin, “A state-of-the-art review of parallel machine scheduling research”, *European Journal of Operational Research*, no. 47, pp. 271–292, 1990.
- [10] B. Chen, C. Potts, and G. Woeginger, “A review of machine scheduling: Complexity, algorithms and approximability”, in *Handbook of Combinatorial Optimization*, pp. 21–169, Kluwer Academic Publishers, 1998.
- [11] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein, *Introduction to Algorithms*. The MIT Press, 3rd ed., 2009.
- [12] W. Rudin, *Principles of Mathematical Analysis*. McGraw-Hill, 3rd ed., 1976.
- [13] H. W. Kuhn, “The hungarian method for the assignment problem”, *Naval Research Logistics Quarterly*, pp. 83–97, 1955.
- [14] A. N. Langville and C. D. Meyer, “Deeper inside pagerank”, *Internet Mathematics*, vol. 1, no. 3, pp. 335–380, 2003.
- [15] G. Gallo, G. Longo, S. Pallottino, and S. Nguyen, “Directed hypergraphs and applications”, *Discrete Applied Mathematics*, vol. 42, pp. 177–201, 1993.



# Appendices



# Appendix A

## Code Listing

We first provide the code for the header files, and then the code for the source ones.

Listing A.1: types.h

```
1 #include <stdio.h>
2
3 #ifndef types_h
4 #define types_h
5
6 typedef struct {
7     int* elems;    // Elements
8     int card;     // Cardinality
9     int repr;     // Representative
10    int id;        // ID (for sorting)
11 } set_t;
12
13 typedef enum {
14     NOEDGE = 0, // No edge between to vertices
15     LTR = 1,   // Edge from a left-hand vertex to a right-hand one (u->v)
16     RTL = -1,  // Edge from a right-hand vertex to a left-hand one (v->u)
17 } adj_t;
18
19 /* Utility procedure to print an array of set_t */
20 static inline void print_sets(set_t* sets, int dim) {
21     for(int i=0;i<dim;i++) {
22         printf("{ ");
23         for(int j=0;j<sets[i].card;j++)
24             printf("%i ", sets[i].elems[j]+1);
25         printf("} [id=%i,card=%i]\n",sets[i].id,sets[i].card);
26     }
27     printf("\n");
28 }
29
30 /* Utility procedure to print one set_t */
31 static inline void print_set(set_t set) {
32     printf("{ ");
33     for(int i=0;i<set.card;i++)
34         printf("%i ", set.elems[i]+1);
35     printf("}");
36 }
37
38 /* Utility functions used for counting_sort */
39 static inline int cardinality(set_t set) {
40     return set.card;
41 }
42
43 static inline int id(set_t set) {
44     return set.id-1; // -1 is necessary for the counting sort to work
45 }
46
47 static inline int representative(set_t set) {
48     return (set.repr != -1) ? set.repr : 0;
49 }
50
51 #endif /* types_h */
```

Listing A.2: rng.h

```

1 #ifndef rng_h
2 #define rng_h
3
4 #define RNG_UNIF 0
5 #define RNG_GAUSS 1
6 #define RNG_EXPO 2
7
8 #define STD_DEV 3.5 // Standard Deviation for the Normal Distribution
9
10 /* lambda for the exponential distribution */
11 int my_rand(int type, int min, int max, double lambda);
12
13 #endif /* rng_h */

```

Listing A.3: queue.h

```

1 #ifndef queue_h
2 #define queue_h
3
4 #include <stdlib.h>
5
6 /* Integer Queue */
7
8 typedef struct queue_elem {
9     int elem;
10    struct queue_elem *next;
11 } queue_elem_t;
12
13 typedef struct {
14     queue_elem_t *head;
15     queue_elem_t *tail;
16     int count;
17 } queue_t;
18
19 static inline queue_t* new_queue() {
20     queue_t *q = malloc(sizeof(*q));
21     if(q == NULL) return NULL;
22
23     q->head = NULL;
24     q->tail = NULL;
25     q->count = 0;
26
27     return q;
28 }
29
30 static inline void enqueue(queue_t* q, int elem) {
31     //printf("enqueue %i\n", elem);
32     if(q == NULL) return;
33
34     queue_elem_t *new = malloc(sizeof(queue_elem_t));
35     if(new == NULL) return;
36
37     new->elem = elem;
38     new->next = NULL;
39
40     if(q->head == NULL) { // Empty queue
41         q->head = new;
42         q->tail = new;
43     }
44     else {
45         q->tail->next = new;
46         q->tail = new;
47     }
48     (q->count)++;
49 }
50
51 static inline int dequeue(queue_t* q) {
52     //printf("Dequeue ");
53     if(q == NULL) return -1;
54
55     if(q->head == NULL) return -1; // Empty queue
56
57     queue_elem_t* old_head = q->head;
58     int to_return = q->head->elem;
59
60     q->head = q->head->next;
61     free(old_head); // Free memory of the old node
62     (q->count)--;
63     //printf("%i\n", to_return);

```

```

64
65 return to_return;
66 }
67
68 static inline void destroy_queue(queue_t* q) {
69     while(q->head != NULL) {
70         queue_elem_t* old = q->head;
71         q->head = q->head->next;
72         free(old);
73     }
74     free(q);
75 }
76 #endif /* queue_h */

```

Listing A.4: intlist.h

```

1 #ifndef intlist_h
2 #define intlist_h
3
4 #include <stdio.h>
5 #include <stdlib.h>
6
7 typedef struct elem {
8     int val;
9     struct elem * next;
10 } elem_t;
11
12 typedef struct {
13     elem_t * head;
14     int length;
15 } intlist_t;
16
17 static inline intlist_t* new_list() {
18     intlist_t* new = malloc(sizeof(intlist_t)); // Placeholder
19     if(new == NULL) return NULL;
20     new->head = NULL;
21     new->length = 0;
22     return new;
23 }
24
25 static inline int extract_head(intlist_t *list) {
26     if(list == NULL || list->head == NULL) return -1;
27     elem_t *tmp = list->head;
28     int to_return = tmp->val;
29
30     list->head = list->head->next;
31     free(tmp);
32
33     (list->length)--;
34
35     return to_return;
36 }
37
38 static inline void insert(intlist_t *list, int elem) {
39     if(list == NULL) return;
40
41     elem_t *new = malloc(sizeof *new);
42     if(new == NULL) return;
43     new->val = elem;
44     new->next = NULL;
45
46     if(list->head == NULL) {
47         list->head = new;
48     }
49     else {
50         elem_t *curr = list->head;
51         while(curr->next != NULL) {
52             curr = curr->next;
53         }
54         curr->next = new;
55     }
56
57     (list->length)++;
58 }
59
60 static inline int is_empty(intlist_t *list) {
61     return (list->head == NULL) ? 1 : 0;
62 }
63

```

```

64 static inline void remove_elem(intlist_t *list, int elem) {
65     if(list == NULL || list->head == NULL) return;
66     if(list->head->next == NULL && list->head->val == elem) {
67         free(list->head);
68         list->head = NULL;
69         return;
70     }
71     elem_t *prev = NULL, *curr = list->head;
72     while(curr != NULL && curr->val != elem) {
73         prev = curr;
74         curr = curr->next;
75     }
76     if(curr != NULL) {
77         if(prev == NULL) {
78             elem_t *tmp = list->head;
79             list->head = list->head->next;
80             free(tmp);
81         }
82         else {
83             prev->next = curr->next;
84             free(curr);
85         }
86     }
87     (list->length)--;
88 }
89 }
90
91 static inline int in_list(intlist_t *list, int elem) {
92     if(list == NULL || list->head == NULL) return 0;
93     elem_t* curr = list->head;
94     while(curr != NULL) {
95         if(curr->val == elem)
96             return 1;
97         curr = curr->next;
98     }
99     return 0;
100 }
101
102 static inline void destroy_list(intlist_t *list) {
103     while(list->head != NULL) {
104         elem_t* tmp = list->head;
105         list->head = list->head->next;
106         free(tmp);
107     }
108     free(list);
109 }
110
111 #endif /* intlíst_h */

```

Listing A.5: countingsort.h

```

1 #ifndef countingsort_h
2 #define countingsort_h
3
4 #include <stdio.h>
5
6 #include "types.h"
7
8 set_t* counting_sort(set_t* sets, int dim, int k, int (*key)(set_t));
9
10 int* integer_counting_sort(int* a, int dim, int max);
11
12 #endif /* countingsort_h */

```

Listing A.6: sm.h

```

1 /**
2  * Algorithms for SEMI-MATCHINGS
3  *
4  */
5
6 #ifndef sm_h
7 #define sm_h
8
9 #include "types.h"
10
11 typedef enum {
12     WHITE,

```

```

13  GREY,
14  BLACK
15 } color_t;
16
17 void asm1(set_t* sets, int m, int n);
18
19 set_t* asm2(set_t** sets, int m, int n);
20
21 #endif /* sm_h */

```

Listing A.7: pagerank.h

```

1 #ifndef pagerank_h
2 #define pagerank_h
3
4 #include "types.h"
5
6 set_t* pagerank_heuristic(set_t* sets, int m, int n, int k); // K = sum of the cardinality
7
8 #endif /* pagerank_h */

```

Listing A.8: avlnode.h

```

1 #ifndef avlnode_h
2 #define avlnode_h
3
4 #include "rng.h"
5
6 /* Nodes of an AVL */
7 typedef struct avl_node_s {
8     struct avl_node_s *left;
9     struct avl_node_s *right;
10    struct avl_node_s *parent;
11    double key;
12    struct node_data_s {
13        int size; // Size and last_position of a
14        int *h; // Hash
15        int *a; // Array
16    } data;
17 } avl_node_t;
18
19 /* Gets a random element from the array a of the field "data" */
20 static inline int avl_node_get_random_element(avl_node_t *node) {
21     return node->data.a[my_rand(RNG_UNIF, 0, node->data.size-1, 1)];
22 }
23
24 /* Adds an element to a node */
25 static inline void avl_node_insert_element(avl_node_t *node, int elem) {
26     node->data.a[node->data.size] = elem;
27     node->data.h[elem] = node->data.size;
28     (node->data.size)++;
29 }
30
31 /* Tells if elem is in node */
32 static inline int avl_node_contains(avl_node_t* node, int elem) {
33     return (node->data.h[elem] != -1) ? 1 : 0;
34 }
35
36 /* Removes an element from a node */
37 static inline void avl_node_remove_element(avl_node_t *node, int elem) {
38     int d = node->data.a[(node->data.size)-1]; // Size-1 = last non-empty position
39     int i = node->data.h[elem]; // Index in the array of the element to be removed
40     node->data.a[i] = d; // The last element of the array
41     node->data.h[d] = i; // Update the index of element d
42     node->data.h[elem] = -1; // Remove the element form the hash table
43     (node->data.size)--;
44 }
45
46
47 #endif /* avlnode_h */

```

Listing A.9: avl.h

```

1 #ifndef avl_h
2 #define avl_h
3
4 #include "types.h"

```

```

5 #include "avlnode.h"
6
7 /* AVL */
8 typedef struct avl_tree_s {
9     avl_node_t *root;
10 } avl_tree_t;
11
12
13 /* Creates a new AVL tree. */
14 avl_tree_t *avl_create(void);
15
16 /* Insert a new node. */
17 void avl_insert(avl_tree_t *tree, double key, int element, int arraysize);
18
19 /* Find the node containing the given key */
20 avl_node_t *avl_find(avl_tree_t *tree, double key);
21
22 /* Finds the first node (in decreasing order) with data.size > 0 */
23 avl_node_t *avl_find_nonempty_maximum(avl_tree_t *tree);
24
25 /* Destroys an AVL */
26 void avl_destroy(avl_tree_t *tree);
27
28 #endif /* avl_h */

```

Listing A.10: greedy.h

```

1 #ifndef greedy_h
2 #define greedy_h
3
4 #include "types.h"
5
6 set_t* pure_greedy_heuristic(set_t** sets, int m, int n);
7
8 set_t* forward_greedy_heuristic(set_t** sets, int m, int n);
9
10 set_t* asm2_heuristic(set_t** sets, int m, int n);
11
12 #endif /* greedy_h */

```

Listing A.11: rng.c

```

1 #include <math.h>
2 #include <float.h>
3 #include <stdlib.h>
4
5 #include "rng.h"
6
7 /* Generatore di numeri casuali normal-distributed */
8 static double generateGaussianNoise(double mu, double sigma) {
9     const double epsilon = DBL_MIN;
10    const double two_pi = 2.0*3.14159265358979323846;
11
12    static double z0, z1;
13    static int generate = 0;
14    generate = !generate;
15
16    if (!generate)
17        return z1 * sigma + mu;
18
19    double u1, u2;
20    do
21    {
22        u1 = rand() * (1.0 / RAND_MAX);
23        u2 = rand() * (1.0 / RAND_MAX);
24    }
25    while ( u1 <= epsilon );
26
27    z0 = sqrt(-2.0 * log(u1)) * cos(two_pi * u2);
28    z1 = sqrt(-2.0 * log(u1)) * sin(two_pi * u2);
29    return z0 * sigma + mu;
30 }
31
32 int my_rand(int type, int min, int max, double lambda) {
33     switch(type) {
34         case RNG_UNIF: {
35             double gen = rand() / (RAND_MAX + 1.0);
36             gen *= (max - min + 1);

```

```

37     gen += min;
38     return (int) gen;
39 }
40 case RNG_GAUSS: {
41     int gen;
42     while((gen = (int) generateGaussianNoise(min + (max - min) / 2.0, (max - min) / 2.0 /
43         STD_DEV)) > max || gen < min);
44     return gen;
45 }
46 case RNG_EXPO: {
47     double gen = rand() / (RAND_MAX + 1.0);
48     double num = (-log(gen) / lambda);
49     while(num < min || num > max) {
50         gen = rand() / (RAND_MAX + 1.0);
51         num = (-log(gen) / lambda);
52     }
53     return (int) floor(num);
54 }
55 default:
56     return 0; // Not yet implemented
57 }

```

Listing A.12: countingsort.c

```

1 #include <stdlib.h>
2
3 #include "countingsort.h"
4
5 set_t* counting_sort(set_t* sets, int dim, int k, int (*key)(set_t)) {
6     int i;
7     int *C = calloc(k+1, sizeof(int));
8     if(C == NULL) return NULL;
9     set_t *B = malloc(dim * sizeof *B);
10    if(B == NULL) return NULL;
11
12    for(i=0;i<dim;i++) {
13        C[key(sets[i])]++;
14    }
15
16    for(i=1;i<=k;i++) {
17        C[i] = C[i-1] + C[i];
18    }
19
20    for(i=dim-1;i>=0;i--) {
21        B[C[key(sets[i])]-1] = sets[i]; // -1 because B starts from 0
22        C[key(sets[i])--];
23    }
24
25    free(C);
26
27    return B;
28 }
29
30 int* integer_counting_sort(int* a, int dim, int max) {
31     int i;
32     int *C = calloc(max+1, sizeof *C);
33     if(C == NULL) return NULL;
34     int *B = malloc(dim * sizeof *B);
35     if(B == NULL) return NULL;
36
37     for(i=0;i<dim;i++) {
38         C[a[i]]++;
39     }
40
41     for(i=1;i<=max;i++) {
42         C[i] = C[i-1] + C[i];
43     }
44
45     for(i=dim-1;i>=0;i--) {
46         B[C[a[i]]-1] = a[i];
47         C[a[i]]--;
48     }
49
50     free(C);
51
52     return B;
53
54 }

```

Listing A.13: asm1.c

```

1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <limits.h>
4
5 #include "sm.h"
6 #include "queue.h"
7 #include "intlist.h"
8 #include "countingsort.h"
9
10 /* BFS on the graph adj, from root, with BF-Tree built in parents */
11 static int bfs(adj_t** adj, int root, int* parents, color_t* colors, int m, int n) {
12     int i;
13     intlist_t* neighbors;
14     queue_t *q = new_queue();
15
16     for(i=0;i<m+n;i++)
17         colors[i] = WHITE;
18     enqueue(q, root);
19     colors[root] = GREY;
20
21     int bestV = -1;
22     int degM_bestV = INT_MAX;
23
24     while(q->count != 0) {
25         int w = (dequeue(q));
26
27         neighbors = new_list();
28
29         if(w < m) { // It's a U vertex (left-hand)
30             for(i=0;i<n;i++) {
31                 if(adj[w][i] == LTR) { // i is an unmatched neighbor
32                     insert(neighbors, i+m); // V-vertices are denoted by their matrix-index + m
33                 }
34             }
35         }
36         else { // It's a V vertex (right-hand)
37             int degM_w = 0;
38             for(i=0;i<m;i++) {
39                 if(adj[i][w-m] == RTL) { // i is a matched neighbor
40                     insert(neighbors, i); // U-vertices are denoted just by their matrix-index
41                     degM_w++;
42                 }
43             }
44             if(bestV == -1 || degM_w < degM_bestV) {
45                 bestV = w-m;
46                 // Update the degM_bestV
47                 degM_bestV = 0;
48                 for(i=0;i<m;i++) {
49                     if(adj[i][bestV] == RTL) {
50                         degM_bestV++;
51                     }
52                 }
53             }
54         }
55         int next = extract_head(neighbors);
56         while(next != -1) { // Scan the entire list
57             if(colors[next] == WHITE) {
58                 parents[next] = w;
59                 enqueue(q, next);
60                 colors[next] = GREY;
61             }
62             next = extract_head(neighbors);
63         }
64
65         destroy_list(neighbors);
66         colors[w] = BLACK;
67     }
68
69     destroy_queue(q);
70     return bestV + m; // Return the best V-vertex found from this root
71 }
72
73 /* Asm1 as described in Semi-Matchings paper */
74 void asm1(set_t* sets, int m, int n) {
75     if(m >= 25000) return;
76     int i, j;
77     /* First step: build the graph via its incidence matrix */
78     adj_t** adj = malloc(m * sizeof(adj_t*)); // Allocate the matrix

```

```

79
80 if(adj == NULL) return;
81 for(i=0;i<m;i++) { // Theta(k)
82     adj[i] = calloc(n, sizeof(adj_t)); // Allocate the rows of the matrix - zeroed to not waste time
83     for(j=0;j<sets[i].card;j++) { // Add the edges to the graph based on the sets
84         adj[i][sets[i].elems[j]] = LTR; // Edge from the set to the variable
85     }
86 }
87
88 int* parents = malloc((m+n) * sizeof(int));
89 if(parents == NULL) return;
90 color_t* vert_colors = malloc((m+n) * sizeof *vert_colors);
91 if(vert_colors == NULL) return;
92
93 for(i=0;i<m;i++) { // m BFS
94     int v = bfs(adj, i, parents, vert_colors, m, n); // Will be > m for sure
95     int u = parents[v]; // Will be < m for sure
96     // Add {u, v} to the matching M - done by setting the direction of edge {u, v} to RTL.
97     adj[u][v-m] = RTL;
98     // Switch edges along path from u to v
99     while(u != i) { // u != root
100         v = parents[u]; // > m
101         adj[u][v-m] = LTR; // Remove {u, v} from matching M
102         u = parents[v];
103         adj[u][v-m] = RTL; // Add {u, v} to matching M
104     }
105 }
106
107 for(i=0;i<m;i++) { // Set the final representative for each set
108     for(j=0;j<n;j++) {
109         if(adj[i][j] == RTL) // Should be just one, and it will be the representative
110             sets[i].repr = j; // The VARIABLE-1, not its position
111     }
112 }
113
114 /* Release all the memory */
115 free(parents);
116 free(vert_colors);
117 for(i=0;i<m;i++)
118     free(adj[i]);
119 free(adj);
120 }

```

Listing A.14: asm2.c

```

1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <limits.h>
4
5 #include "sm.h"
6 #include "queue.h"
7 #include "intlist.h"
8 #include "countingsort.h"
9
10 int DFS(adj_t** adj, int vert, int original_root, intlist_t* S, int* load, int m, int n, int*
    parents, int* start) {
11     int i;
12
13     for(i=0;i<m;i++) {
14         if(adj[i][vert] == RTL) { // For each u in MatchedNeighbors(vert) -- u is our i
15             parents[i] = vert+m; // vert is in v
16
17             for(int w=0;w<n;w++) {
18                 if(adj[i][w] == LTR) { // For each w in UnmatchedNeighbors(i)
19                     if(!in_list(S, w) || load[w] > load[vert]) { //degM(adj, w, m, n) > degM(adj, vert, m, n)
20                         { // w already visited or load on w worse than load on v
21                             continue; // skip to next w
22                         }
23                         remove_elem(S, w);
24                         parents[w+m] = i; // Because w is in V
25                         if(load[w] <= load[original_root] - 2) { // From the original root!
26                             *start = w; // Start of the path
27                             return 1;
28                         }
29                     }
30                 }
31             }
32     }

```

```

33 return 0;
34 }
35
36 int find_cost_reducing_path(adj_t** adj, int* load, int m, int n, int* parents, int* root_to_return) {
37     int i;
38     intlist_t *not_visited = new_list();
39     /* At the beginning, not_visited = V */
40     for(i=0;i<n;i++)
41         insert(not_visited, i);
42
43     while(!is_empty(not_visited)) {
44         /* Look for a V-vertex with max degM(v) */
45         int max_degM = -1;
46         int root = -1;
47         elem_t* curr = not_visited->head;
48         while(curr != NULL) {
49             // Calcolo il deg_M del vertice curr
50             int degM_v = load[curr->val]; //degM(adj, curr->val, m, n);
51
52             if(degM_v > max_degM) {
53                 max_degM = degM_v;
54                 root = curr->val;
55             }
56             curr = curr->next;
57         }
58         /* The V-vertex chosen is in root */
59
60         remove_elem(not_visited, root); // Remove root from not_visited
61         parents[root+m] = -1; // Because parents are all together, v indices > m
62         int start_path;
63
64         if(DFS(adj, root, not_visited, load, m, n, parents, &start_path) == 1) {
65             *root_to_return = root; // Path from start_path to root
66             destroy_list(not_visited);
67             return start_path;
68         }
69     }
70     destroy_list(not_visited);
71
72     return -1; // No cost-reducing path exists
73 }
74
75 /* Asm2 as described in Semi-Matchings paper */
76 set_t* asm2(set_t** sets, int m, int n) {
77     int i, j;
78     /* First step: sort the U-vertices (the sets) by increasing degree (cardinality) */
79     set_t* sorted_sets = counting_sort(*sets, m, n, cardinality); // Sort the U vertices by increasing
80     degree
81
82     /* Build the graph via its incidence matrix */
83     adj_t** adj = malloc(m * sizeof(adj_t*)); // Allocate the matrix
84     int* deg = calloc(n, sizeof *deg); // Allocate the deg array that will contain the degree of every
85     V-vertex
86     int* degM = calloc(n, sizeof *deg); // Allocate the degM array that will contain the degree of
87     every V-vertex in the matching M
88
89     if(adj == NULL) return NULL;
90     if(deg == NULL) return NULL;
91     if(degM == NULL) return NULL;
92
93     for(i=0;i<m;i++) { // Theta(k)
94         adj[i] = calloc(n, sizeof(adj_t)); // Allocate the rows of the matrix - zeroed to not waste time
95         for(j=0;j<sorted_sets[i].card;j++) { // Add the edges to the graph based on the sets
96             adj[i][sorted_sets[i].elems[j]] = LTR; // Edge from the set to the variable
97             deg[sorted_sets[i].elems[j]]++; // Count the degree (occurrences) of this variable
98         }
99     }
100
101     /* GREEDY INITIALIZATION */
102     int *load = calloc(n, sizeof *load); // The load is equal to the degree of a V-vertex in the
103     matching M
104     if(load == NULL) return NULL;
105     /* Assign each U-vertex with its least-loaded V-neighbor (in case of ties, choose the V-neighbor
106     with minimum degree) */
107     for(i=0;i<m;i++) {
108         int load_min = INT_MAX;
109         int deg_min = INT_MAX;
110         int v_min = 0;
111         for(j=0;j<sorted_sets[i].card;j++) { // Search the least-loaded V-neighbor
112             if(load[sorted_sets[i].elems[j]] < load_min) {

```

```

107     load_min = load[sorted_sets[i].elems[j]];
108     v_min = sorted_sets[i].elems[j];
109 }
110 }
111
112 for(j=0;j<sorted_sets[i].card;j++) { // Break ties
113     if(load[sorted_sets[i].elems[j]] == load_min && deg[sorted_sets[i].elems[j]] < deg_min) {
114         deg_min = deg[sorted_sets[i].elems[j]];
115         v_min = sorted_sets[i].elems[j];
116     }
117 }
118
119 /* v_min is the least-loaded/least-degree V-neighbor of i */
120 adj[i][v_min] = RTL; // Add {u, v} to matching M
121 load[v_min]++; // And thus increase the load of v_min
122 }
123 /* END: GREEDY INITIALIZATION */
124
125 /* Next: Find and remove cost-reducing paths */
126 int* parents = malloc((m+n) * sizeof(int));
127 if(parents == NULL) return NULL;
128 int root, start;
129
130 while((start = find_cost_reducing_path(adj, load, m, n, parents, &root)) != -1) { // While there
131     // exists a cost-reducing path
132     /* Switch matching and non-matching edges along P */
133     // The cost-reducing path is from start to root
134
135     // Switch edges along cost-reducing path from start to root
136     int v = start; // < m
137     int u = parents[v+m];
138     while(v-m != root) { // u != root (root < m)
139         adj[u][v] = RTL; // Add {v, u} to matching M
140         load[v]++; // Increase the load of this v
141
142         v = parents[u]; // > m
143
144         adj[u][v-m] = LTR; // Remove {u, v} from matching M
145         load[v-m]--; // Decrease the load of this v
146
147         u = parents[v];
148     }
149     for(i=0;i<m+n;i++) // Reset the parents
150         parents[i] = -1;
151 }
152
153 for(i=0;i<m;i++) { // Set the final representative for each set
154     for(j=0;j<n;j++) {
155         if(adj[i][j] == RTL) // Should be just one, and it will be the representative
156             sorted_sets[i].repr = j; // The VARIABLE-1, not its position
157     }
158 }
159
160 set_t* sets_to_return = counting_sort(sorted_sets, m, m+1, id);
161
162 /* Release all the memory */
163 for(i=0;i<m;i++) {
164     free(adj[i]);
165 }
166 free(adj);
167 free(deg);
168 free(degM);
169 free(load);
170 free(parents);
171 free(*sets);
172 free(sorted_sets);
173
174 return sets_to_return;
175 }

```

Listing A.15: avl.c

```

1 #define _XOPEN_SOURCE 500 /* Enable certain library functions (strdup) on linux. See
   feature_test_macros(7) */
2
3 #include <time.h>
4 #include <stdlib.h>
5 #include <stdio.h>

```

```

6 #include <limits.h>
7 #include <string.h>
8 #include <assert.h>
9
10 #include "avl.h"
11
12 /* Creates a new AVL tree. */
13 avl_tree_t *avl_create() {
14     avl_tree_t *tree = NULL;
15
16     if((tree = malloc(sizeof( avl_tree_t ))) == NULL)
17         return NULL;
18
19     tree->root = NULL;
20
21     return tree;
22 }
23
24 /* Initialize a new node. */
25 static avl_node_t *avl_create_node(double key, int arraysize) {
26     avl_node_t *node = NULL;
27
28     if( ( node = malloc( sizeof( avl_node_t ) ) ) == NULL ) {
29         return NULL;
30     }
31
32     node->left = NULL;
33     node->right = NULL;
34     node->parent = NULL;
35     node->key = key;
36     node->data.size = 0;
37     node->data.a = malloc(arraysize * sizeof(int));
38     if(node->data.a == NULL)
39         return NULL;
40     node->data.h = malloc(arraysize * sizeof(int));
41     if(node->data.h == NULL)
42         return NULL;
43
44     return node;
45 }
46
47 /* Find the height of an AVL node recursively */
48 static int avl_node_height( avl_node_t *node ) {
49     int height_left = 0;
50     int height_right = 0;
51
52     if(node->left) height_left = avl_node_height(node->left);
53     if(node->right) height_right = avl_node_height(node->right);
54
55     return height_right > height_left ? ++height_right : ++height_left;
56 }
57
58 /* Find the balance of an AVL node */
59 static int avl_balance_factor( avl_node_t *node ) {
60     int bf = 0;
61
62     if(node->left) bf += avl_node_height(node->left);
63     if(node->right) bf -= avl_node_height(node->right);
64
65     return bf;
66 }
67
68 /* Left Left Rotate */
69 static avl_node_t *avl_rotate_leftleft( avl_node_t *node ) {
70     avl_node_t *a = node;
71     avl_node_t *b = a->left;
72
73     a->left = b->right;
74     if(b->right != NULL) b->right->parent = a; ///////
75     b->right = a;
76     if(b != NULL) b->parent = a->parent; ///////
77     a->parent = b; ///////
78
79
80     return b;
81 }
82
83 /* Left Right Rotate */
84 static avl_node_t *avl_rotate_leftright( avl_node_t *node ) {

```

```

85  avl_node_t *a = node;
86  avl_node_t *b = a->left;
87  avl_node_t *c = b->right;
88
89  a->left = c->right;
90  if(c->right != NULL) c->right->parent = a; //////////
91  b->right = c->left;
92  if(c->left != NULL) c->left->parent = b; //////////
93  c->left = b;
94  b->parent = c; //////////
95  c->right = a;
96  c->parent = a->parent; //////////
97  a->parent = c; //////////
98  return c;
99 }
100
101 /* Right Left Rotate */
102 static avl_node_t *avl_rotate_rightleft( avl_node_t *node ) {
103     avl_node_t *a = node;
104     avl_node_t *b = a->right;
105     avl_node_t *c = b->left;
106
107     //print_albero(a);
108     a->right = c->left;
109     if(c->left != NULL) c->left->parent = a; //////////
110     b->left = c->right;
111     if(c->right != NULL) c->right->parent = b; //////////
112     c->right = b;
113     b->parent = c; //////////
114     c->left = a;
115     c->parent = a->parent; //////////.
116     a->parent = c; //////////
117
118     return c;
119 }
120
121 /* Right Right Rotate */
122 avl_node_t *avl_rotate_rightright( avl_node_t *node ) {
123     avl_node_t *a = node;
124     avl_node_t *b = a->right;
125
126     a->right = b->left;
127     if(b->left != NULL) b->left->parent = a; //////////
128     b->left = a;
129     if(b != NULL) b->parent = a->parent; //////////
130     if(a != NULL) a->parent = b; //////////
131
132     return( b );
133 }
134
135 /* Balance a given node */
136 static avl_node_t *avl_balance_node( avl_node_t *node ) {
137     return node;
138     avl_node_t *newroot = NULL;
139
140     /* Balance our children, if they exist. */
141     if(node->left)
142         node->left = avl_balance_node(node->left);
143     if(node->right)
144         node->right = avl_balance_node(node->right);
145
146     int bf = avl_balance_factor(node);
147
148     if(bf >= 2) {
149         /* Left Heavy */
150
151         if(avl_balance_factor(node->left) <= -1)
152             newroot = avl_rotate_leftright(node);
153         else
154             newroot = avl_rotate_leftleft(node);
155     } else if(bf <= -2) {
156         /* Right Heavy */
157
158         if(avl_balance_factor(node->right) >= 1)
159             newroot = avl_rotate_rightleft(node);
160         else
161             newroot = avl_rotate_rightright(node);
162
163

```

```

164 } else {
165     /* This node is balanced -- no change. */
166     newroot = node;
167 }
168
169 return newroot;
170 }
171
172 /* Balance a given tree */
173 static void avl_balance( avl_tree_t *tree ) {
174     return;
175     avl_node_t *newroot = NULL;
176     newroot = avl_balance_node(tree->root);
177
178     if(newroot != tree->root) {
179         tree->root = newroot;
180         tree->root->parent = NULL;
181     }
182 }
183
184 /* Insert a new node. */
185 void avl_insert( avl_tree_t *tree, double key, int element, int arraysize ) {
186     avl_node_t *node = NULL;
187     avl_node_t *next = NULL;
188     avl_node_t *last = NULL;
189
190     /* Well, there must be a first case */
191     if( tree->root == NULL ) {
192         node = avl_create_node(key, arraysize);
193         avl_node_insert_element(node, element);
194
195         tree->root = node;
196
197         /* Okay. We have a root already. Where do we put this? */
198     } else {
199         next = tree->root;
200
201         while( next != NULL ) {
202             last = next;
203
204             if( key < next->key ) {
205                 next = next->left;
206
207             } else if( key > next->key ) {
208                 next = next->right;
209
210                 /* Have we already inserted this node? */
211             } else if( key == next->key ) {
212                 avl_node_insert_element(next, element);
213                 return; // No need to balance, just append the element in the list!
214             }
215         }
216
217         node = avl_create_node(key, arraysize);
218         avl_node_insert_element(node, element);
219
220         if( key < last->key ) last->left = node;
221         if( key > last->key ) last->right = node;
222
223         node->parent = last;
224     }
225 }
226
227 avl_balance(tree);
228 }
229
230 /* Find the node containing a given key */
231 avl_node_t *avl_find( avl_tree_t *tree, double key ) {
232     avl_node_t *current = tree->root;
233
234     while( current && current->key != key ) {
235         if( key > current->key )
236             current = current->right;
237         else
238             current = current->left;
239     }
240
241     return current;
242 }

```

```

243
244 static avl_node_t *avl_maximum(avl_node_t *node) {
245     if(node == NULL)
246         return NULL;
247
248     while(node->right != NULL)
249         node = node->right;
250
251     return node;
252 }
253
254 static avl_node_t *avl_predecessor(avl_node_t *node) {
255     if(node == NULL)
256         return NULL;
257
258     if(node->left != NULL)
259         return avl_maximum(node->left);
260
261     avl_node_t *parent = node->parent;
262     while(parent != NULL && node == parent->left) {
263         node = parent;
264         parent = parent->parent;
265     }
266
267     return parent;
268 }
269
270 /* Finds the first node (in decreasing order) with data.size > 0 */
271 avl_node_t *avl_find_nonempty_maximum(avl_tree_t *tree) {
272     avl_node_t *maximum = avl_maximum(tree->root);
273     if(maximum->data.size == 0) {
274         maximum = avl_predecessor(maximum);
275         while(maximum != NULL && maximum->data.size == 0) {
276             maximum = avl_predecessor(maximum);
277         }
278     }
279
280     return maximum;
281 }
282
283 /* Frees the memory held by a node */
284 static void destroy_node(avl_node_t *node) {
285     if(node == NULL) return;
286
287     if(node->left)
288         destroy_node(node->left);
289     if(node->right)
290         destroy_node(node->right);
291
292     if(node->data.a != NULL) free(node->data.a);
293     if(node->data.h != NULL) free(node->data.h);
294     free(node);
295 }
296
297 /* Destroys the entire tree */
298 void avl_destroy(avl_tree_t *tree) {
299     destroy_node(tree->root);
300     free(tree);
301 }

```

Listing A.16: greedy.c

```

1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <limits.h>
4 #include <math.h>
5
6 #include "greedy.h"
7 #include "countingsort.h"
8
9 set_t* pure_greedy_heuristic(set_t** sets, int m, int n) { // Runs in O(k) time
10     int i, j;
11     int candidate;
12     int min_freq;
13     //print_sets(sets, m);
14     set_t* sorted_sets = counting_sort(*sets, m, n, cardinality); // O(m+n)
15     int* frequencies = calloc(n, sizeof(int));
16     if(frequencies == NULL) return NULL;
17

```

```

18  for(i=0;i<m;i++) { // O(k)
19      min_freq = INT_MAX;
20      candidate = -1;
21      for(j=0;j<sorted_sets[i].card;j++) {
22          if(frequencies[sorted_sets[i].elems[j]] < min_freq) {
23              min_freq = frequencies[sorted_sets[i].elems[j]]; // Select the least-chosen variable from
                each set, using direct access
24              candidate = sorted_sets[i].elems[j];
25          }
26      }
27      sorted_sets[i].repr = candidate;
28      frequencies[candidate]++;
29  }
30
31  set_t* sets_to_return = counting_sort(sorted_sets, m, m, id); // O(m + n) ---- IDs go from 1 to m+1
32
33  free(frequencies);
34  free(sorted_sets);
35  free(*sets);
36  return sets_to_return;
37 }
38
39 /* Select the least-appearing var, if it does not exceed the mean.
40 * If it does, choose according the pure greedy strategy (i.e. the least-chosen)
41 */
42 set_t* forward_greedy_heuristic(set_t** sets, int m, int n) {
43     int i, j;
44     //print_sets(sets, m);
45     set_t* sorted_sets = counting_sort(*sets, m, n, cardinality);
46     int* occurrences = calloc(n, sizeof(int));
47     if(occurrences == NULL) return NULL;
48     int* frequencies = calloc(n, sizeof(int));
49     if(frequencies == NULL) return NULL;
50     int mean = (int) ceil((double) m/n);
51     int candidate_freq, candidate_occ, min_freq, min_occ;
52
53     for(i=0;i<m;i++)
54         for(j=0;j<sorted_sets[i].card;j++)
55             occurrences[sorted_sets[i].elems[j]]++; // Count the occurrences of all the variables
56
57     for(i=0;i<m;i++) {
58         candidate_freq = -1;
59         candidate_occ = -1;
60         min_freq = INT_MAX;
61         min_occ = INT_MAX;
62         for(j=0;j<sorted_sets[i].card;j++) {
63             if(occurrences[sorted_sets[i].elems[j]] < min_occ && frequencies[sorted_sets[i].elems[j]] <
                mean) {
64                 candidate_occ = sorted_sets[i].elems[j];
65                 min_occ = occurrences[sorted_sets[i].elems[j]];
66             }
67             if(frequencies[sorted_sets[i].elems[j]] < min_freq) {
68                 candidate_freq = sorted_sets[i].elems[j];
69                 min_freq = frequencies[sorted_sets[i].elems[j]];
70             }
71         }
72         if(candidate_occ == -1) // Did not find a future variable chosen less than the mean time
73             sorted_sets[i].repr = candidate_freq; // Pure greedy
74         else
75             sorted_sets[i].repr = candidate_occ; // candidate_occ is the least-appearing var that does not
                exceed the mean (>=)
76
77         /* Keep the occurrences updated */
78
79         frequencies[sorted_sets[i].repr]++; // Increase the frequency of the chosen variable
80         for(j=0;j<sorted_sets[i].card;j++)
81             occurrences[sorted_sets[i].elems[j]]--; // Decrease the occurrences of the variables in this set
82     }
83
84     set_t* sets_to_return = counting_sort(sorted_sets, m, m, id); // O(m + n) ---- IDs go from 1 to m+1
85
86     free(frequencies);
87     free(occurrences);
88     free(sorted_sets);
89     free(*sets);
90     return sets_to_return;
91 }
92
93 set_t* asm_2_heuristic(set_t** sets, int m, int n) {

```

```

94     int i, j;
95     int candidate;
96     int min_freq, min_occ;
97     set_t* sorted_sets = counting_sort(*sets, m, n, cardinality); // O(m+n)
98     int* frequencies = calloc(n, sizeof(int));
99     if(frequencies == NULL) return NULL;
100    int* occurrences = calloc(n, sizeof(int));
101    if(occurrences == NULL) return NULL;
102
103    for(i=0;i<m;i++)
104        for(j=0;j<sorted_sets[i].card;j++)
105            occurrences[sorted_sets[i].elems[j]]++; // Count the occurrences of all the variables
            (the degrees)
106
107    for(i=0;i<m;i++) { // O(k)
108        min_freq = INT_MAX;
109        min_occ = INT_MAX;
110        candidate = -1;
111        for(j=0;j<sorted_sets[i].card;j++) {
112            if(frequencies[sorted_sets[i].elems[j]] < min_freq) {
113                min_freq = frequencies[sorted_sets[i].elems[j]]; // Select the least-chosen variable
                from each set, using direct access
114                candidate = sorted_sets[i].elems[j];
115            }
116        }
117        for(j=0;j<sorted_sets[i].card;j++) {
118            /* Tie-breaking mechanism */
119            if(frequencies[sorted_sets[i].elems[j]] == min_freq &&
120                occurrences[sorted_sets[i].elems[j]] < min_occ) {
121                min_occ = occurrences[sorted_sets[i].elems[j]];
122                candidate = sorted_sets[i].elems[j];
123            }
124            /* This heuristic does not decrease the occurrences of the variables in the current set */
125        }
126        sorted_sets[i].repr = candidate;
127        frequencies[candidate]++;
128    }
129
130    set_t* sets_to_return = counting_sort(sorted_sets, m, m, id); // O(m + n) ---- IDs go from 1 to
        m+1
131
132    free(frequencies);
133    free(occurrences);
134    free(sorted_sets);
135    free(*sets);
136    return sets_to_return;
137 }

```

Listing A.17: pagerank.c

```

1 #include <stdlib.h>
2 #include <float.h>
3 #include <math.h>
4
5 #include "pagerank.h"
6 #include "avlnode.h"
7 #include "avl.h"
8 #include "intlist.h"
9 #include "rng.h"
10
11 /* Inside the sets there are the variables -1! */
12
13 void print_matrix(double **W, int m, int n) {
14     for(int i=0;i<m;i++) {
15         for(int j=0;j<n;j++) {
16             printf("%lf\t",W[i][j]);
17         }
18         printf("\n");
19     }
20 }
21
22 set_t* pagerank_heuristic(set_t* sets, int m, int n, int k) { // K = sum of the cardinality
23     int i, j;
24     /* Weights */
25     double **W = malloc(m * sizeof *W);
26     double *total_weights = calloc(n, sizeof *total_weights);
27     int *ones = calloc(n, sizeof *ones);
28     int *not_zero = calloc(n, sizeof *not_zero);

```

```

29 double max_total_weight = 0.;
30 int var_max_total_weight = -1;
31 /* AVLS */
32 avl_tree_t **avls = malloc(n * sizeof *avls);
33
34 if(W == NULL) return NULL;
35 if(total_weights == NULL) return NULL;
36 if(avls == NULL) return NULL;
37 if(ones == NULL) return NULL;
38 if(not_zero == NULL) return NULL;
39
40 /* Initialize the AVLS */
41 for(i=0;i<n;i++)
42     avls[i] = avl_create(); /* They are indexed by the variables */
43
44 /* Initialize the Weights matrix and the variables' AVLS */
45 for(i=0;i<m;i++) { // Theta(k)
46     W[i] = calloc(n + 1, sizeof *(W[i])); // n+1 to store the current cardinality in the last column
47     for(j=0;j<sets[i].card;j++) {
48         int var = sets[i].elems[j];
49         W[i][var] = (double) 1./sets[i].card; // Populate the matrix...
50         total_weights[var] += (double) 1./sets[i].card; // ... and update the total weights of the
51         variables
52         if(total_weights[var] > max_total_weight) { // Update the max total weight
53             max_total_weight = total_weights[var];
54             var_max_total_weight = var;
55         }
56         not_zero[var]++;
57         /* The variable var appears in the i-th set, with relative weight of W[i][var] */
58         if(W[i][var] != 1.) {
59             //printf("Inserisco nell'avl %lf della variabile %i\n", W[i][var], var);
60             avl_insert(avls[var], W[i][var], i, m); // Insert in the avls[var]: the key W[i][var] which
61             is the relative weight of var in the i-th set
62         }
63         else { // Singleton
64             ones[var]++;
65             sets[i].repr = var;
66         }
67     }
68     W[i][n] = sets[i].card;
69 }
70
71 /* Core of the algorithm */
72 for(i=0;i<k-m;i++) {
73     /* Pick the maximum total weight -> max_total_weight */
74     /* Where is the var of maximum total weight with its heighest relative weight */
75
76     avl_node_t *maximum_node = avl_find_nonempty_maximum(avls[var_max_total_weight]); // Mi da il
77     primo nodo (massimo) in cui A.size != 0
78
79     int chosen_set = avl_node_get_random_element(maximum_node); // Dammi un insieme con quella
80     variabile e pi alto relative weight
81
82     avl_node_remove_element(maximum_node, chosen_set); // Rimuovi dal nodo dell'avl quell'insieme
83     (chosen_set)
84
85     /* Zero the variable and decrease the total maximum */
86     total_weights[var_max_total_weight] -= W[chosen_set][var_max_total_weight];
87     W[chosen_set][var_max_total_weight] = 0.;
88     not_zero[var_max_total_weight]--; // One more zero added to this set
89     W[chosen_set][n]--; // Decrease the logical cardinality of chosen_set
90
91     max_total_weight = 0;
92     for(j=0;j<n;j++) { // Update all the variables in this set
93         if(W[chosen_set][j] != 0.) {
94             avl_node_remove_element(avl_find(avls[j], W[chosen_set][j]), chosen_set); // This variable
95             does not appear in this subset with this set anymore
96
97             total_weights[j] -= W[chosen_set][j]; // Update total weights
98             W[chosen_set][j] = (double) 1. / W[chosen_set][n]; // 1 / current cardinality
99
100             if(W[chosen_set][j] == 1) { // Do we have a new representative?
101                 sets[chosen_set].repr = j;
102                 ones[j]++;
103             }
104
105             total_weights[j] += W[chosen_set][j];
106             if(W[chosen_set][j] != 1.) { // Only re-add this variable if it is not the only left in the set
107                 avl_insert(avls[j], W[chosen_set][j], chosen_set, m); // Add back the new relative weight

```

```

102     }
103     }
104
105     if(total_weights[j] >= max_total_weight && not_zero[j] != ones[j]) { // Update the max total
weight
106         max_total_weight = total_weights[j];
107         var_max_total_weight = j;
108     }
109 }
110 }
111
112 /* Release the memory */
113 free(total_weights);
114 for(i=0;i<n;i++)
115     avl_destroy(avls[i]);
116 free(avls);
117 for(i=0;i<m;i++)
118     free(W[i]);
119 free(W);
120 free(ones);
121
122 return sets;
123 }

```

Listing A.18: main.c

```

1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <math.h>
4 #include <dirent.h>
5 #include <unistd.h>
6 #include <time.h>
7 #include <float.h>
8 #include <sys/types.h>
9 #include <sys/stat.h>
10
11 #include "types.h"
12 #include "sm.h"
13 #include "greedy.h"
14 #include "pagerank.h"
15 #include "countingsort.h"
16 #include "rng.h"
17
18 #define timing(f) start = clock(); \
19 (f); \
20 end = clock(); \
21 time_spent = (double) (end - start) / CLOCKS_PER_SEC;
22
23 typedef struct {
24     double* time_asm1;
25     double* time_asm2;
26     double* time_greedy;
27     double* time_greedyfwd;
28     double* time_asm2heur;
29     double* time_pagerank;
30     int errors_greedy;
31     int errors_greedyfwd;
32     int errors_pagerank;
33     int errors_asm2heur;
34     int badness_greedy;
35     int badness_greedyfwd;
36     int badness_asm2heur;
37     int badness_pagerank;
38 } alg_stats_t;
39
40 static int is_regular_file(const char *path)
41 {
42     struct stat path_stat;
43     stat(path, &path_stat);
44     return S_ISREG(path_stat.st_mode);
45 }
46
47 static set_t* build_sets_from_userinput(int *m, int *n, int **occ) {
48     set_t* sets;
49
50     printf("How many sets? ");
51     scanf("%i", m);
52     printf("How many vars? ");
53     scanf("%i", n);

```

```

54
55 *occ = calloc(*n, sizeof *occ);
56 if(*occ == NULL) return NULL;
57 sets = malloc(*m * sizeof(*sets));
58 if(sets == NULL) return NULL;
59
60 for(int i=0;i<*m;i++) {
61     printf("How many elements in set %i? ", i+1);
62     scanf("%i", &sets[i].card);
63     sets[i].elems = malloc(sets[i].card * sizeof(int));
64     if(sets[i].elems == NULL) return NULL;
65     sets[i].id = i+1; // For the counting sort to work
66     sets[i].repr = -1;
67
68     for(int j=0;j<sets[i].card;j++) {
69         scanf("%i", &sets[i].elems[j]);
70         sets[i].elems[j]--;
71         ((*occ)[sets[i].elems[j]])++;
72     }
73 }
74
75 return sets;
76 }
77
78 static set_t* build_sets_from_file(const char *filename, int *m, int *n, int **occ) {
79     set_t* sets;
80
81     if(is_regular_file(filename)) { // File
82         FILE* fileinput = fopen(filename, "r");
83         if(fileinput == NULL) return NULL;
84
85         fscanf(fileinput, "%i", m);
86         fscanf(fileinput, "%i", n);
87
88         *occ = calloc(*n, sizeof *occ);
89         if(*occ == NULL) return NULL;
90         sets = malloc(*m * sizeof *sets);
91         if(sets == NULL) return NULL;
92
93         for(int i=0;i<*m;i++) {
94             /* Format: CARDINALITY {x1 x2 x3 ... xn}\n */
95             fscanf(fileinput, "%i {", &sets[i].card);
96             sets[i].id = i+1; // For the counting sort to work
97             sets[i].repr = -1;
98             sets[i].elems = malloc(sets[i].card * sizeof(*(sets[i].elems)));
99
100             for(int j=0;j<sets[i].card;j++) {
101                 fscanf(fileinput, "%i", &sets[i].elems[j]);
102                 sets[i].elems[j]--;
103                 ((*occ)[sets[i].elems[j]])++;
104             }
105
106             fscanf(fileinput, "}\n");
107         }
108
109         fclose(fileinput);
110         return sets;
111     }
112     else { // Directory
113         /* Placeholder */
114     }
115     return NULL;
116 }
117
118 /* Build M sets with N vars, rand_type_vars distributed. Cardinality distribution type is
119    rand_type_card */
120 static set_t* build_sets_randomly(int m, int n, int **occ, int rand_type_card, int rand_type_vars) {
121     set_t* sets = malloc(m * sizeof *sets);
122     if(sets == NULL) return NULL;
123     *occ = calloc(n, sizeof *occ);
124     if(*occ == NULL) return NULL;
125     int k, ok;
126
127     for(int i=0;i<m;i++) {
128         sets[i].card = my_rand(rand_type_card, 1, n, (double) 1./((double) n/2.));
129         sets[i].elems = malloc(sets[i].card * sizeof(int));
130         if(sets[i].elems == NULL) return NULL;
131         sets[i].id = i+1; // For the counting sort to work
132         sets[i].repr = -1;

```

```

132
133     for(int j=0;j<sets[i].card;j++) {
134         ok = 1;
135         sets[i].elems[j] = my_rand(rand_type_vars, 1, n, (double) 1./n);
136         for(k=0;k<j;k++) {
137             if(sets[i].elems[k] == sets[i].elems[j]-1) { // Already generated this var, try another one
138                 j--;
139                 ok = 0;
140                 break;
141             }
142         }
143         if(ok) { // Successfully added a new variable
144             sets[i].elems[j]--;
145             ((*occ)[sets[i].elems[j]])++;
146         }
147     }
148
149     int* tmp = integer_counting_sort(sets[i].elems, sets[i].card, n); // Sort the variables in the set
150     free(sets[i].elems);
151     sets[i].elems = tmp;
152 }
153
154 return sets;
155 }
156
157 static int check_max_freq(set_t* sets, int m, int n, int* frequencies, int* var_of_max_freq) {
158     int i;
159     for(i=0;i<n;i++)
160         frequencies[i] = 0;
161
162     for(i=0;i<m;i++) { // Print the solution and count the frequencies
163         frequencies[sets[i].repr]++;
164         //printf("%i from ", sets[i].repr+1);
165         //print_set(sets[i]);
166         //printf("\n");
167     }
168     /* Search the maximum frequency */
169     int maximum_frequency = 0;
170     *var_of_max_freq = -1;
171     for(i=0;i<n;i++) {
172         if(frequencies[i] > maximum_frequency) {
173             maximum_frequency = frequencies[i];
174             *var_of_max_freq = i;
175         }
176     }
177
178     return maximum_frequency;
179 }
180
181 static set_t* simplify(set_t* sets, int* m, int* n, int* occ) {
182     int i, j;
183     int mean = (int) ceil((double) (*m)/(*n));
184     int discarded = 0;
185
186     for(i=0;i<*m;i++) {
187         for(j=0;j<sets[i].card;j++) {
188             if(occ[sets[i].elems[j]] < mean) {
189                 sets[i].repr = sets[i].elems[j];
190                 discarded++;
191                 printf("[SIMPLIFY] Took %i from ", sets[i].repr+1);
192                 print_set(sets[i]);
193                 printf("\n");
194                 break;
195             }
196         }
197     }
198
199     if(discarded > 0) {
200         set_t* sorted_sets = counting_sort(sets, *m, *n, representative);
201
202         for(i=(*m)-discarded;i<*m;i++)
203             free(sorted_sets[i].elems);
204
205         sorted_sets = realloc(sorted_sets, ((*m)-discarded) * sizeof *sorted_sets);
206
207         *m = (*m) - discarded;
208
209         sorted_sets = counting_sort(sorted_sets, *m, (*m) + discarded, id);
210

```

```

211     for(i=0;i<*m;i++)
212         sorted_sets[i].id = i+1;
213
214     return sorted_sets;
215 }
216 return sets;
217 }
218
219
220 void call_algorithms(set_t** sets, int m, int n, int *input_size, alg_stats_t *stats) {
221     int* frequencies;
222     int maximum_frequency;
223     int var_of_max_freq;
224     int i;
225     static int last_time_position = 0;
226     static int sum_k = 0;
227     int optimum;
228
229     /* Variables used for timing */
230     clock_t start, end;
231     double time_spent;
232
233     /* Prepare the frequencies array */
234     frequencies = calloc(n, sizeof *frequencies);
235     if(frequencies == NULL) return;
236
237     /* Call Asm1 */
238     printf("#### Asm1 ####\n");
239     timing(asm1(*sets, m, n)); // Call Asm1
240     if(m < 25000) {
241         maximum_frequency = check_max_freq(*sets, m, n, frequencies, &var_of_max_freq); // Count the
242             maximum frequency
243         printf("The maximum frequency is %i (variable %i). Time: %lf sec.\n", maximum_frequency,
244             var_of_max_freq+1, time_spent);
245         optimum = maximum_frequency;
246     }
247     else {
248         optimum = -1;
249     }
250     if(stats != NULL && stats->time_asm1 != NULL)
251         stats->time_asm1[last_time_position] = time_spent; // Store the time it took
252
253     /* Call Asm2 */
254     printf("\n#### Asm2 ####\n");
255     for(i=0;i<m;i++)
256         (*sets)[i].repr = -1; // Remove the previous representatives
257     timing(*sets = asm2(sets, m, n));
258
259     maximum_frequency = check_max_freq(*sets, m, n, frequencies, &var_of_max_freq); // Count the
260         maximum frequency
261     printf("The maximum frequency is %i (variable %i). Time: %lf sec.\n", maximum_frequency,
262         var_of_max_freq+1, time_spent);
263     if(stats != NULL && stats->time_asm2 != NULL)
264         stats->time_asm2[last_time_position] = time_spent; // Store the time it took
265
266     if(optimum == -1) optimum = maximum_frequency; // Remember the result of Asm2 as the optimal value
267
268     /* Call the pure greedy heuristic */
269     printf("\n#### PURE GREEDY HEURISTIC ####\n");
270     for(i=0;i<m;i++)
271         (*sets)[i].repr = -1; // Remove the previous representatives
272     timing(*sets = pure_greedy_heuristic(sets, m, n));
273
274     maximum_frequency = check_max_freq(*sets, m, n, frequencies, &var_of_max_freq); // Count the
275         maximum frequency
276     printf("The maximum frequency is %i (variable %i). Time: %lf sec.\n", maximum_frequency,
277         var_of_max_freq+1, time_spent);
278     if(stats != NULL && stats->time_greedy != NULL)
279         stats->time_greedy[last_time_position] = time_spent; // Store the time it took
280     if(maximum_frequency != optimum) {
281         stats->badness_greedy += maximum_frequency-optimum;
282         (stats->errors_greedy)++; // Check if the result is the same as the one of Asm1
283     }
284
285     /* Call the forward greedy heuristic */
286     printf("\n#### FORWARD GREEDY HEURISTIC ####\n");
287     for(i=0;i<m;i++)
288         (*sets)[i].repr = -1; // Remove the previous representatives
289     timing(*sets = forward_greedy_heuristic(sets, m, n));

```

```

284
285 maximum_frequency = check_max_freq(*sets, m, n, frequencies, &var_of_max_freq); // Count the
    maximum frequency
286 printf("The maximum frequency is %i (variable %i). Time: %lf sec.\n", maximum_frequency,
    var_of_max_freq+1, time_spent);
287 if(stats->time_greedyfwd != NULL)
288     stats->time_greedyfwd[last_time_position] = time_spent; // Store the time it took
289 if(maximum_frequency != optimum) {
290     stats->badness_greedyfwd += maximum_frequency-optimum;
291     (stats->errors_greedyfwd)++; // Check if the result is the same as the one of Asm1
292 }
293
294 /* Call the Asm2 heuristic */
295 printf("\n### ASM2 INIT HEURISTIC ###\n");
296 for(i=0;i<m;i++)
297     (*sets)[i].repr = -1; // Remove the previous representatives
298 timing(*sets = asm2_heuristic(sets, m, n));
299
300 maximum_frequency = check_max_freq(*sets, m, n, frequencies, &var_of_max_freq); // Count the
    maximum frequency
301 printf("The maximum frequency is %i (variable %i). Time: %lf sec.\n", maximum_frequency,
    var_of_max_freq+1, time_spent);
302 if(stats->time_asm2heur != NULL)
303     stats->time_asm2heur[last_time_position] = time_spent; // Store the time it took
304 if(maximum_frequency != optimum) {
305     stats->badness_asm2heur += maximum_frequency-optimum;
306     (stats->errors_asm2heur)++; // Check if the result is the same as the one of Asm1
307 }
308
309 /* Call the Pagerank heuristic */
310 printf("\n### PAGERANK HEURISTIC ###\n");
311 int k = 0; // Count the cardinalities of all the sets
312 for(i=0;i<m;i++) {
313     (*sets)[i].repr = -1; // Remove the previous representatives
314     k += (*sets)[i].card;
315 }
316 sum_k += k;
317 timing(*sets = pagerank_heuristic(*sets, m, n, k));
318
319 maximum_frequency = check_max_freq(*sets, m, n, frequencies, &var_of_max_freq); // Count the
    maximum frequency
320 printf("The maximum frequency is %i (variable %i). Time: %lf sec.\n", maximum_frequency,
    var_of_max_freq+1, time_spent);
321 if(stats->time_pagerank != NULL)
322     stats->time_pagerank[last_time_position] = time_spent; // Store the time it took
323 if(maximum_frequency != optimum) {
324     stats->badness_pagerank += maximum_frequency-optimum;
325     (stats->errors_pagerank)++; // Check if the result is the same as the one of Asm1
326 }
327
328 if(stats->time_asm1 != NULL || stats->time_asm2 != NULL || stats->time_greedy != NULL ||
    stats->time_greedyfwd != NULL || stats->time_pagerank != NULL)
329     last_time_position++; // Update the index of the time arrays
330
331 if(input_size != NULL)
332     *input_size = sum_k;
333
334 free(frequencies);
335 }
336
337 int main(int argc, const char * argv[]) {
338     int m, n;
339     set_t* sets;
340     int i;
341     int *occurrences;
342
343     /* Set the random seed */
344     unsigned int seed = (argc == 5) ? (unsigned int) atoi(argv[4]) : (unsigned int) time(NULL); // Save
        the seed for debug purposes;
345     printf("RNG Seed: %u\n\n", seed);
346     srand(seed);
347
348     /* Build the input instance */
349     switch(argc) {
350     case 5:
351     case 4: { // m sets, n vars, x different instances
352         int instances = atoi(argv[3]);
353         n = atoi(argv[2]);
354         m = atoi(argv[1]);

```

```

355 FILE* stats_file = fopen("stats_nuove_unif.csv", "a");
356 if(stats_file == NULL) return -3;
357 /* Create the arrays for storing execution times of the algorithms */
358 alg_stats_t stats = {
359     malloc(instances * sizeof(double)),
360     malloc(instances * sizeof(double)),
361     malloc(instances * sizeof(double)),
362     malloc(instances * sizeof(double)),
363     malloc(instances * sizeof(double)),
364     malloc(instances * sizeof(double)),
365     0,
366     0,
367     0,
368     0,
369     0,
370     0,
371     0,
372     0
373 };
374 double *t_asm1 = stats.time_asm1;
375 double *t_asm2 = stats.time_asm2;
376 double *t_greedy = stats.time_greedy;
377 double *t_greedyfwd = stats.time_greedyfwd;
378 double *t_pagerank = stats.time_pagerank;
379 double *t_asm2heur = stats.time_asm2heur;
380 double min_time_asm1 = DBL_MAX, min_time_asm2 = DBL_MAX, min_time_greedy = DBL_MAX,
min_time_greedyfwd = DBL_MAX, min_time_pagerank = DBL_MAX, min_time_asm2heur = DBL_MAX;
381 double max_time_asm1 = -1., max_time_asm2 = -1., max_time_greedy = -1., max_time_greedyfwd =
-1., max_time_pagerank = -1., max_time_asm2heur = -1.;
382 double sum_asm1 = 0., sum_asm2 = 0., sum_greedy = 0., sum_greedyfwd = 0., sum_pagerank = 0.,
sum_asm2heur = 0.;
383 int total_input_size = 0;
384 int simplify_sets = 0;
385 int removed_vars = 0;
386
387 clock_t start = clock();
388
389 for(i=0;i<instances;i++) {
390     /* Build and print the input */
391     sets = build_sets_randomly(m, n, &occurrences, RNG_UNIF, RNG_GAUSS);
392     printf("[%i] m = %i, n = %i\n", i, m, n);
393     simplify_sets += m; // For statistical purposes
394     //print_sets(sets, m);
395     sets = simplify(sets, &m, &n, occurrences);
396     printf("\nAfter Simplifying: m = %i, n = %i\n", m, n);
397     simplify_sets -= m;
398
399     /* Check which variables have been simplified */
400     for(int j=0;j<n;j++)
401         occurrences[j] = 0;
402
403     for(int x=0;x<m;x++) {
404         for(int j=0;j<sets[x].card;j++) {
405             occurrences[sets[x].elems[j]]++;
406         }
407     }
408
409     for(int j=0;j<n;j++) {
410         if(occurrences[j] == 0)
411             removed_vars++;
412     }
413
414     //print_sets(sets, m);
415
416     /* Call the algorithms on this instance */
417     call_algorithms(&sets, m, n, &total_input_size, &stats);
418
419     /* Update the time stats */
420     if(t_asm1[i] < min_time_asm1)
421         min_time_asm1 = t_asm1[i]; // Min
422     if(t_asm1[i] > max_time_asm1)
423         max_time_asm1 = t_asm1[i]; // Max
424     sum_asm1 += t_asm1[i]; // Mean
425     /* Time stats for ASM2 */
426     if(t_asm2[i] < min_time_asm2)
427         min_time_asm2 = t_asm2[i]; // Min
428     if(t_asm2[i] > max_time_asm2)
429         max_time_asm2 = t_asm2[i]; // Max
430     sum_asm2 += t_asm2[i]; // Mean

```



```

499         (stats.errors_pagerank > 0) ? (double) stats.badness_pagerank / stats.errors_pagerank : 0.,
500         (stats.errors_asm2heur > 0) ? (double) stats.badness_asm2heur /
stats.errors_asm2heur : 0.,
501         (int) ceil((double) simplify_sets / instances),
502         (int) ceil((double) removed_vars / instances),
503         sum_asm1 / (double) instances,
504         sum_asm2 / (double) instances,
505         sum_greedy / (double) instances,
506         sum_greedyfwd / (double) instances,
507         sum_pagerank / (double) instances,
508         min_time_asm1,
509         min_time_asm2,
510         min_time_greedy,
511         min_time_greedyfwd,
512         min_time_pagerank,
513         max_time_asm1,
514         max_time_asm2,
515         max_time_greedy,
516         max_time_greedyfwd,
517         max_time_pagerank
518     );
519
520     fclose(stats_file);
521     free(t_asm1);
522     free(t_asm2);
523     free(t_greedy);
524     free(t_greedyfwd);
525     free(t_pagerank);
526     free(t_asm2heur);
527     return 0; // quit the entire program
528 }
529 case 2:
530     sets = build_sets_from_file(argv[1], &m, &n, &occurrences);
531     break;
532 default:
533     sets = build_sets_from_userinput(&m, &n, &occurrences);
534     break;
535 }
536
537 if(sets == NULL) {
538     fprintf(stderr, "There was an error during the creation of the instance. [QUIT]\n");
539     return -1;
540 }
541
542 printf("Input: m = %i, n = %i\n", m, n);
543 print_sets(sets, m);
544 sets = simplify(sets, &m, &n, occurrences);
545 printf("\nAfter Simplifying: m = %i, n = %i\n", m, n);
546 print_sets(sets, m);
547
548 alg_stats_t stats = {
549     NULL,
550     NULL,
551     NULL,
552     NULL,
553     NULL,
554     NULL,
555     0,
556     0,
557     0,
558     0,
559     0,
560     0,
561     0,
562     0
563 };
564 call_algorithms(&sets, m, n, NULL, &stats);
565 /* Free the sets */
566 for(i=0; i<m; i++)
567     free(sets[i].elems);
568 free(sets);
569 free(occurrences);
570 return 0;
571 }

```