

1 An Efficient Approximate Greedy Algorithm to Solve 2 the Minimum Set Cover Problem

3 Jorge Delgado, Héctor Ferrada and Cristóbal A. Navarro

4 *Instituto de Informática, Universidad Austral de Chile, Chile*

5 Abstract

The Minimum Set Cover Problem (MSCP) is a combinatorial optimization problem belonging to the NP-Hard class in computer science. For this reason, there is no algorithm that in the worst case ensures finding an optimal solution in polynomial time. For a given universe \mathcal{X} , the popular greedy heuristic, called **Greedy-SetCover**, is the main theoretical contribution to obtain an approximate solution for the MSCP in polynomial time, offering an optimal approximate ratio $(1 - o(|\mathcal{X}|)) \ln |\mathcal{X}|$. In this article, we propose an approximate algorithm for MSCP within a succinct representation of the input dataset, whose empirical performance improves **Greedy-SetCover** both in quality and execution time, while offering the same optimal approximation ratio for the problem. Our experiments show that the proposed algorithm is magnitudes of times faster than the aforementioned greedy one, obtaining on average a cardinality much closer to the optimal solution. Furthermore, because we work on a succinct representation that allows us to compute operations between sets using bitwise operators, we can process much larger datasets than state-of-the-art solutions. As a result, our proposal also is a suitable alternative for processing large datasets as required by the current Big Data era.

6 *Keywords:* Combinatorial optimization, Set cover, NP-Hard, Greedy
7 algorithms, Approximate algorithms

8 1. Introduction

9 A *set-covering* for a given universe \mathcal{X} is a group of subsets of \mathcal{X} whose
10 union includes all the elements of \mathcal{X} . Since there can be many different combinations of sets whose union gives the same universe, in some circumstances
11

12 it is desired to select only the best combination, or optimal, among all the
 13 available sets. When we define the optimum as the smallest number of sets to
 14 cover \mathcal{X} and we require to determine what this optimal solution is, then we
 15 are talking about the Minimum Set Cover Problem (MSCP). The MSCP is
 16 an NP-hard optimization problem [17, 22] that we formally define as follows.
 17 Given an instance $(\mathcal{X}, \mathcal{F})$ as input dataset, where $\mathcal{F} = \{S_1, \dots, S_m\}$ is the
 18 family of subsets and \mathcal{X} is the universe formed by the union of all sets in \mathcal{F} ,
 19 the problem consists in to find a set-covering $\mathcal{C} \subseteq \mathcal{F}$ of \mathcal{X} ; i.e., $\mathcal{X} = \bigcup_{S \in \mathcal{C}} S$ of
 20 minimum size. According to the definition it is possible that there is more
 21 than one solution for MSCP, in which case all of them are called optimal so-
 22 lutions for the given instance $(\mathcal{X}, \mathcal{F})$, and to solve MSCP, it is enough to find
 23 any of them. The Figure 1 shows an universe of twelve elements distributed
 24 in seven subsets $S_i \in \mathcal{F}, 1 \leq i \leq |\mathcal{F}|$ which are the input instance $(\mathcal{X}, \mathcal{F})$
 25 for the problem, and for this input dataset the unique optimal solution is
 26 $\mathcal{C}^* = \{S_1, S_4, S_5\}$.

27

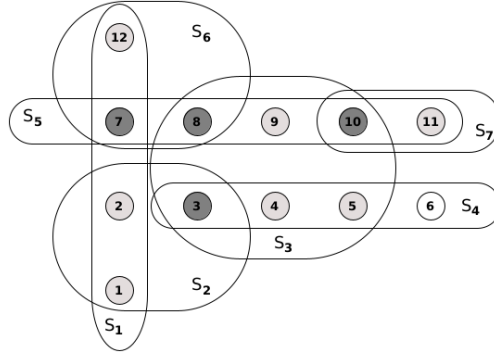


Figure 1: Illustration of an instance $(\mathcal{X}, \mathcal{F})$ of the MSCP, with $|\mathcal{X}| = 12$ elements and $\mathcal{F} = \{S_1, \dots, S_7\}$. The optimal solution for this example is the minimum set-covering $\mathcal{C}^* = \{S_1, S_4, S_5\}$.

28 The MSCP is widely used in many optimization problems that require
 29 resources to be allocated. As a small example in combinatorics, within a
 30 multidisciplinary research project, we may need to select the minimum num-
 31 ber of researchers to cover all the research areas and skills \mathcal{X} necessary to
 32 carry out the project. Thus, a solution for MSCP it is enough. MSCP also
 33 appears naturally in tasks of coordination and distribution of finite resources
 34 within the industry. For instance, within industrial planning, in order to

35 avoid overlaps and/or reduce operating costs, an optimal solution may be
 36 required to allocate worker shifts. In the same line of planning, the Transit
 37 Route Network Design Problem (TrNDP) can be derived from the MSCP
 38 [29] and use this solution to establish transport routes of a company. In an-
 39 other application, some model-based view planning methods also require a
 40 set-covering to solve the view planning problem (VPP), which involves deter-
 41 mining the minimum number of views that must be used to reconstruct the
 42 part in 3D [32]. In industrial robot coordination, a solution for MSCP also
 43 is used to minimize synchronization points required by its algorithms [34]. A
 44 set-covering also appears as an optimization problem to make an efficient dis-
 45 tribution of tasks in an *Internet of Things* (IoT) network within the context
 46 of limited energy [11]. In *Crew scheduling*, the problem of assigning a group
 47 of workers (a *crew*) to a set of tasks can be formulated from the MSCP[4].
 48 Contardo et al., [16], propose a hybrid algorithm that iterates between com-
 49 puting a set-covering and constructing a Voronoi diagram to solve MSCP
 50 in geometric contexts. The MSCP has many other applications in different
 51 areas of research on optimization problems, such as industrial applications
 52 [13], information retrieval [35], in graph theory to determine a transversal in
 53 hypergraphs [14], in satellite systems for traffic allocation [33], among others.

54
 55 In the worst case, any current method that determines the optimal com-
 56 bination of sets from \mathcal{F} , requires at least exponential time in the size of the
 57 family of sets, $|\mathcal{F}| = m$. We can homologate all the possible combinations
 58 from \mathcal{F} of an exhaustive algorithm as if it were a binary counter of m bits,
 59 starting from 1 to 2^m . The binary counter determines each possible sets com-
 60 bination, where if the i th bit is one in a combination ρ , $1 \leq i \leq m$, then the
 61 set S_i is included in ρ . Thus, we need to evaluate each one of the 2^m states of
 62 the counter, checking if the combinations correspond to a true set-covering
 63 for \mathcal{X} , and remember the configuration with the fewest sets. In order to check
 64 if a group of sets includes all the elements of \mathcal{X} , it is necessary to perform
 65 a validation on all $n = |\mathcal{X}|$ elements, which requires $\Omega(n)$ time by using a
 66 sequential method that evaluates each element one by one. Therefore, the
 67 worst case time of an exhaustive optimal algorithm is $\Omega(n \cdot 2^m)$.

68 In this sense, there are many heuristics to improve, in some cases, the
 69 execution time of the exhaustive search that obtains the optimal solution
 70 to the problem. For example, we can narrow the search space by determin-
 71 ing the smallest number of sets that any optimal solution must contain, as
 72 follows. If we find the smallest integer k such that the sum of the sizes of

the k largest sets of \mathcal{F} reaches the size of the universe, then we know that any optimal solution must contain at least k subsets. For that, we need to know what are the sizes of the k largest sets in \mathcal{F} , which we can obtain by ordering the family of subsets in $O(|\mathcal{F}| \lg |\mathcal{F}|)$ time. Thus, having found this k , we start the search by evaluating any combination of k sets; if we do not find a set-covering for the universe, then we iterate again increasing k by one unit and repeating the search for this updated k value. When we find a set-covering for \mathcal{X} , we stop the iteration reporting the last k value as the size of the optimal solution.

As we have seen, the MSCP has a great impact on academia, scientific research and several industry fields. However, under the assumption that it is very unlikely to find a polynomial-time algorithm to obtain the optimal set-covering, because its NP-completeness, great part of the effort of the scientific community is in the development of polynomial-time approximations for the problem. Meaning to offer faster solutions at the cost of a lower quality. Furthermore, we believe that in combination with efficient programming techniques, it is possible to build faster and/or smaller solutions to the problem, either to determine the optimal configuration or to design new approximate solutions.

The rest of this article is structured as follows, Section 2 details the related work in the field, the next Section 3 presents the new greedy solution for MSCP. Section 4 shows experimental results and finally Section 5 presents conclusions with future work.

2. Related Work

In 1972 Karp showed that the MSCP is NP-Hard [22], and therefore, very unlikely to be optimally solved in polynomial time. Two decades later, Lund and Yannakakis based on the work of Karp, showed in [27] that it is not possible to obtain a polynomial time approximate algorithm for solving the MSCP, with ratio better than $\frac{1}{4} \lg n$, unless NP has slightly superpolynomial time deterministic algorithms, that is $O(n^{\text{poly} \log n})$, where $n = |\mathcal{X}|$ is the universe size. Thus, every polynomial solution must have an approximate ratio $\geq \frac{1}{4} \log n$, showing that it is also difficult that a polynomial approximation method becomes $\leq \frac{1}{2} \log_2 n \simeq 0.72 \ln n$. A few years later, in 1998, Uriel Feige [18] adjusted the approximation ratio to $(1 - o(n)) \ln n$, but always con-

109 sidering that NP does not include slightly superpolynomial time algorithms.
 110 Therefore, given the NP-completeness of the MSCP which indicates that ev-
 111 ery optimal solution has superpolynomial time in the worst case, most of the
 112 current optimal strategies focus their success on the development of heuristics
 113 to prune the search space. The latter has the risk that for certain cases of the
 114 input instance it is not possible to do any pruning, or narrow the search only
 115 a little. For this reason, in the literature we find several optimal algorithms
 116 based on different pruning strategies for the search tree [2, 3, 5]. Among this
 117 type of algorithms, the method proposed by Lan et al. [25] finds the optimal
 118 set-covering by applying a search strategy called meta-heuristics for random
 119 priority search, introducing randomness into the resolution of the problem.
 120 Other exhaustive search algorithms are mainly based on branch-and-bound
 121 and branch-and-cut methods [1, 6, 19].

122
 123 Maybe the major efforts in MSCP research are the design of new approx-
 124 imations, avoiding the restrictions imposed by the Hardness of the problem.
 125 In this line, there are several proposals in the state-of-the-art that try to
 126 obtain faster solutions and desirably close to the optimum quality. In 1973
 127 Johnson published his Greedy algorithm [21] which determines an approx-
 128 imation for the MSCP in polynomial time. The algorithm 1, included in
 129 [17], describes the pseudocode of Johnson’s method. His intuitive heuristic
 130 consists of iteratively selecting the set that covers the most number of re-
 131 maining uncovered elements in the universe, subtracting the included items
 132 from the current universe and iterating until all items are covered. In 1979
 133 Chvatal [15] investigated the work of Johnson to propose an approximate
 134 solution with very similar characteristics, which can be reduced to the John-
 135 son’s results under certain conditions. On the other hand, an interesting
 136 result appeared in [10], where the authors identified two issues that arise
 137 when determining a set-covering with metaheuristic approaches: solution in-
 138 feasibility and set redundancy¹. Accordingly, they propose a solution on a
 139 1-flip neighbourhood structure to find a set-covering with metaheuristic by
 140 including a penalty approach to obtain an approximate solution to the prob-
 141 lem. Lim et al., [26] also analyzed the main greedy algorithm approach to

¹An incomplete set-covering is considered to be infeasible if one or more of the elements
 of the universe are uncovered (i.e., not really a solution). A set is considered to be
 redundant if all the elements covered by the set are also covered by other sets in the
 solution.

the MSCP by evaluating eager and lazy versions of the algorithm. Caprara et al.[12] proposed an heuristic based on Lagrangian relaxation to find an approximation for the problem, whose implementation won an important competition of that decade, called FASTER. Among randomized algorithms we highlight the method of Peuzin-Jubert et al., in [31], which introduces an heuristic called randomized rounding algorithm for MSCP, achieving an approximation ratio of $k(1 - (c/n)^{1/k})$ for some small constant c .

Approximate minimum set-covering can also be treated with parallelized methods, for example, Koufogiannakis and Young [24] present parallel and distributed δ -approximation algorithms to cover problems, such as weighted set cover or weighted vertex cover, where δ is the maximum number of variables on which any constraint depends.

Gomes et al. [20] carried out an empirical study of approximation algorithms for the Vertex Cover and the Set Covering Problems. They tested the classical greedy algorithm against some variants of it, showing that, in general, the greedy strategy [21] is the fastest and the best approximation for the MSCP.

2.1. The Greedy Set Cover Algorithm

The solution with major acceptance in the literature is the greedy approximation introduced by Johnson in 1973 [21], which is also described in Cormen et al.[17]. This method is a polynomial time algorithm whose heuristic selects at each iteration the most promising subset among those available. The key idea is to select the set $S \in \mathcal{F}$ that contains the largest number of elements that have not yet been included by any set of the partial solution that has been built up to that iteration. The process ends the iterations when all the elements of the universe have been considered by at least one of the selected sets. This strategy ensures finding a solution \mathcal{C} whose approximation relation is $\frac{|\mathcal{C}|}{|\mathcal{C}^*|} \leq \ln |\mathcal{X}| + 1$, where \mathcal{C}^* is an optimal solution for the problem. The pseudocode described in 1 corresponds to the greedy set covering algorithm detailed in Cormen et al.[17] with $O(\min\{|\mathcal{X}|, |\mathcal{F}|\} \cdot |\mathcal{X}| \cdot |\mathcal{F}|)$ complexity time, in the worst case. This theoretical time is explained by the loop of lines 4...8, where it is necessary to determine the set S that maximizes $|\mathcal{U} \cap S|$. To do that, in each one of the $\min\{|\mathcal{X}|, |\mathcal{F}|\}$ iterations, we may need to examine up to $|\mathcal{F}|$ subsets of size $|\mathcal{X}|$, resulting in a bound of $O(|\mathcal{X}| \cdot |\mathcal{F}|)$ time.

Algorithm 1 GREEDY-SETCOVER algorithm in [17] to compute a set covering for the universe \mathcal{X} and the family of subsets \mathcal{F}

Require: The universe \mathcal{X} and the family of subsets \mathcal{F} whose union is \mathcal{X}

```

1: procedure GREEDY-SETCOVER( $\mathcal{X}, \mathcal{F}$ )
2:    $\mathcal{U} = \mathcal{X}$ 
3:    $\mathcal{C} = \emptyset$ 
4:   while  $\mathcal{U} \neq \emptyset$  do                                      $\triangleright O(\min\{|\mathcal{X}|, |\mathcal{F}|\})$  time
5:     select  $S \in \mathcal{F}$  that maximizes  $|\mathcal{U} \cap S|$               $\triangleright O(|\mathcal{U}| \cdot |\mathcal{F}|)$  time
6:      $\mathcal{U} = \mathcal{U} - S$ 
7:      $\mathcal{C} = \mathcal{C} \cup S$ 
8:   end while
9:   return  $\mathcal{C}$ 
10: end procedure

```

179 The correctness of Algorithm 1 is ensured by the loop condition, which does
180 not allow to finish the process until all the elements of the universe have
181 been considered in the solution. The quality of the greedy algorithm has an
182 approximation ratio of $O(\ln |\mathcal{X}|)$ [17], which is theoretically optimal consid-
183 ering the conditions for the problem discussed in the Introduction (Section
184 1).

185 3. An Approximate Greedy Set Cover Algorithm in a Succinct 186 Data Representation

187 The classic greedy strategy described above achieves its time and quality
188 by selecting the set with the largest number of uncovered elements by the
189 partial solution. This polynomial-time algorithm 1 is essentially the best
190 theoretical research result to date regarding to approximate solutions for the
191 set cover problem. However, the method can take too long when the size of
192 the data collection is large. Mainly because iteratively performs exhaustive
193 searches for the best set among all available sets in \mathcal{F} , validating in the worst
194 case up to $|\mathcal{F}|$ sets. Furthermore, the exaggerated generality of the algorithm
195 invites the development of more adequate data structures to represent the
196 sets; for instance, by using an environment that facilitates the computation
197 of set operations such as Union and Intersection. These aspects have opened
198 up a research opportunity that has given rise to this contribution.
199 Derived from Algorithm 1, we have designed a set cover approximated algo-
200 rithm incorporating new heuristics to empirically improve the original greedy

method. Thus achieving greater speedup and space savings, as well as guaranteeing optimal asymptotic complexity in terms of approximation quality —considering that this problem is encapsulated by the NP-hard class. We highlight two key points in our proposal: *i*- the design of an efficient and lightweight data structure that allows calculations of operations between sets directly in RAM; and *ii*- a new heuristic to prune the space search of the best set in each iteration. The final part of our proposal is an implementation of the new algorithm that improves both the execution time of the original method and the quality of the approximation.

When analyzing Algorithm 1, we can notice that the key step is in line 5, when it selects the next set $S \in \mathcal{F}$ to be included in the solution. This selection has to check all the remaining sets of \mathcal{F} , obligating to execute an exhaustive search in \mathcal{F} , requiring $O(|\mathcal{X}| \cdot |\mathcal{F}|)$ time in the worst case. The time of this selection can become very significant for large sizes of the universe \mathcal{X} or the family of subsets \mathcal{F} . This is the most expensive task of each iteration in the algorithm, because most of these exhaustive searches are performed on a large subset of \mathcal{F} .

Consider Figure 1, where the optimal solution is $\{S_1, S_4, S_5\}$. The classic greedy algorithm starts by selecting the subset S_3 to build the solution $\mathcal{C} = \{S_3, S_1, S_4, S_7\}$, or including S_5 instead of S_7 in \mathcal{C} , whose size $|\mathcal{C}| = 4$. However, note that we can select S_4 as the first subset of the solution \mathcal{C} , because the element number 6, the rightmost in S_4 , is not included by any other set. In this way, we can reduce the search space if we focus only on covering the elements that are present in the smallest number of subsets. Also hoping that, with this way of selecting each set $S \in \mathcal{F}$, other elements of the universe will also be covered, resulting in a possible faster and more efficient algorithm than 1 and also with a better approximation ratio. Thus, we iteratively select the set containing the largest number of uncovered elements that are present in the smallest number of subsets in \mathcal{F} . Formally, if an item $x \in \mathcal{X}$ is included exactly in k sets, then we call x as an *element of grade k* . In the example of the Figure 1, the white element, the item with identifier 6, has grade $k = 1$; the light gray elements (identifiers 1, 2, 4, 5, 9, 11, 12) have grade $k = 2$; and the dark gray elements (identifier 3, 7, 8, 10) have grade $k = 3$. We will then use this notion in order to provide a new greedy algorithm to compute the set-covering for the instance $(\mathcal{X}, \mathcal{F})$. Bellow we formally detail our algorithm, which pseudocode is described in 2.

Algorithm 2 SUCCINCT-SETCOVER algorithm to compute a set cover of the universe \mathcal{X} and the family of subsets \mathcal{F} .

Require: The universe of elements \mathcal{X} and the family of subsets \mathcal{F} whose union is \mathcal{X}

```

1: procedure SUCCINCT-SETCOVER( $\mathcal{X}, \mathcal{F}$ )
2:   be MAP a map data structure of length  $|\mathcal{X}|$  to store, for each  $x_i \in \mathcal{X}$ , the lists of
   ID-subsets where  $x_i$  is.
3:   MAP  $\leftarrow$  CREATEMAP( $\mathcal{X}, \mathcal{F}$ )  $\triangleright O(|\mathcal{X}| \cdot |\mathcal{F}|)$  time
4:   SORTUNIVERSE(MAP,  $\mathcal{F}$ )  $\triangleright O(|\mathcal{X}| \log |\mathcal{X}|)$  time
5:    $\mathcal{C} = \emptyset$ 
6:    $\mathcal{U} = \mathcal{X}$ 
7:    $ini = k = 1$ 
8:   while  $\mathcal{U} \neq \emptyset$  do  $\triangleright O(\min\{|\mathcal{X}|, |\mathcal{F}|\})$  time
9:      $end \leftarrow$  FINDRANGE(MAP,  $ini, k$ )  $\triangleright O(|\mathcal{X}|)$  time for all the calls
10:     $\ell = end - ini + 1$   $\triangleright$  numbers of elements of grade  $k$  to cover
11:    while  $\ell > 0$  do
12:       $\langle S, \delta \rangle \leftarrow$  FINDSUBSET(MAP,  $ini, end$ )  $\triangleright O(\frac{|\mathcal{X}|}{W} |\mathcal{F}|)$  time
13:       $\mathcal{U} = \mathcal{U} - S$   $\triangleright O(\frac{|\mathcal{X}|}{W})$  time
14:       $\mathcal{C} = \mathcal{C} \cup S$   $\triangleright O(\frac{|\mathcal{X}|}{W})$  time
15:       $\ell = \ell - \delta$ 
16:    end while
17:     $ini = end + 1$ 
18:     $k = k + 1$ 
19:  end while
20:  return  $\mathcal{C}$ 
21: end procedure

```

240 We first detail the function of each procedure included in the pseudocode
 241 of the algorithm.

- 242 • CREATEMAP(\mathcal{X}, \mathcal{F}), procedure that creates and returns a map data
 243 structure, called MAP, of pairs $(x_i, L_i[])$, where for each $x_i \in \mathcal{X}$, $L_i[]$ is
 244 the list of all subsets $S_i \in \mathcal{F}$ in which x_i is. It can be built in $O(|\mathcal{X}| \cdot |\mathcal{F}|)$
 245 time by simple iterative inspection.
- 246 • SORTUNIVERSE(MAP, \mathcal{F}), procedure that sorts the MAP structure of
 247 size $|\mathcal{X}|$ in increasing order with respect to the number of times that
 248 each $x_i \in \mathcal{X}$ appears in the subsets of \mathcal{F} . In this process, for each x_i ,
 249 with the same MAP, we obtain in constant time the length of each list
 250 of subsets that contain x_i . Thus, MAP is sorted in $O(|\mathcal{X}| \log |\mathcal{X}|)$ time.

- 251 • FINDRANGE(MAP, ini, k), procedure that determines and returns the
 252 right limit end of the range MAP[$ini \dots end$] of all elements $x \in \mathcal{X}$ of
 253 grade k . Thanks that MAP is sorted, determining each interval $I =$
 254 MAP[$ini \dots end$] can be done in linear time over the size $|I|$. If there
 255 are no elements of grade k , then returns $ini - 1$.
- 256 • FINDSUBSET(MAP, ini, end), procedure that determines and returns
 257 the pair $\langle S, \delta \rangle$, where $S \in \mathcal{F}$ is the subset that contains the major num-
 258 ber δ of minimal uncovered elements inside the range MAP[$ini \dots end$]
 259 of all elements with the same grade. If there are ties, then select
 260 the subset that maximizes $|S \cap \mathcal{U}|$ among all the lists MAP[x_i].list,
 261 $\forall x_i \in \text{MAP}[ini \dots end]$.

262 As input, we are given the instance $(\mathcal{X}, \mathcal{F})$, which are the universe of
 263 elements \mathcal{X} and the family of subsets \mathcal{F} , whose union is \mathcal{X} . We also require
 264 to build a MAP data structure² to associate to each $x \in \mathcal{X}$ the list of all
 265 subsets $S \in \mathcal{F}$ where x is —created at the line 3 in the pseudocode. At line
 266 4, we sort the universe \mathcal{X} in MAP (by the length of the list that contains each
 267 element) by calling to the procedure SORTUNIVERSE. We build the solution
 268 \mathcal{C} incrementally, which is initially an empty set —see line 5. We start iter-
 269 ating with $\mathcal{U} = \mathcal{X}$ while \mathcal{U} has elements, which means that the size of our
 270 solution \mathcal{C} is lower than $|\mathcal{X}|$ —line 8. The process finds, at line 9, each range
 271 MAP[$ini \dots end$] of all elements $x \in \mathcal{X}$ of grade k , starting with $k = ini = 1$.
 272 These elements of grade k , are also called *Minimal Elements*, because we look
 273 for uncovered elements of minimal grade k among all the remaining elements
 274 not yet included by the partial solution \mathcal{C} . Therefore, at the beginning, we
 275 start searching in MAP the interval of all elements that are included only by
 276 $k = 1$ subsets of \mathcal{F} ; that is, in the first group of iterations (the inner loop at
 277 lines 11—16) the minimal elements are items of degree $k = 1$. Thus, in the
 278 next iterations of the while-cycle of lines 8—19 we search for the next interval
 279 of minimal items of grade $k = 2$. As in the MAP structure, the ranges of
 280 elements are grouped by their degree value, placing one after the other as
 281 the degree of the elements increases. Therefore, we now loop in the inner
 282 while-cycle, lines 11—16, while the range size ℓ be greater than 0. It means

²This corresponds to a BST or a Red/Black Tree [17, 23], where each key has associate an additional data structure such as `std::map`, which is provided by the GNU C++ Standard Library and implements a Red/Black Tree.

283 that there are uncovered elements of degree k and, consequently, we must
 284 continue including subsets in the solution until we cover all the minimal ele-
 285 ments of the range.

286

287 The correctness of our algorithm is given by the condition of the outer
 288 loop of line 8, which does not end as long as there is at least one uncovered
 289 element.

290 With the example of Figure 1, unlike the classic greedy algorithm that
 291 starts by searching exhaustively among all the subsets in \mathcal{F} to select the
 292 subset S_3 , the first search of the SUCCINCT-SETCOVER algorithm is only
 293 among the sets that contains at least one element of minimal degree, since
 294 we start with $k = 1$, the only set that satisfies the condition is S_4 . The
 295 next selected set, for $k = 2$, is S_1 , which once included elements of degree 2
 296 remain uncovered, so S_5 is chosen in the second iteration of the inner loop
 297 for $k = 2$. Note that in the example we never search for elements of degree
 298 $k = 3$ because these were implicitly included when we added the sets to cover
 299 the elements of degree $k = 1$ and $k = 2$.

300 3.1. Asymptotic Analysis of the SUCCINCT-SETCOVER Algorithm

301 We detail the asymptotic analysis of our SUCCINCT-SETCOVER Algo-
 302 rithm 2 for both, the time consumption and the quality of the solution.

303

304 To obtain a threshold of the time, let us break down the pseudocode
 305 of the algorithm. At the beginning we need to build a sorted MAP dictio-
 306 nary, requiring a time complexity of $O(|\mathcal{X}| \cdot |\mathcal{F}| + |\mathcal{X}| \log |\mathcal{X}|) = O(|\mathcal{X}| \cdot |\mathcal{F}|)$,
 307 as described above. Later, at line 9, all the calls for the procedure FIND-
 308 RANGE require $O(|\mathcal{X}|)$ time, which is achieved by taking advantage of the
 309 order in MAP, performing a simple scan. At line 12, each call of the pro-
 310 cedure FINDSUBSET(MAP, ini , end) requires a time of $len \cdot O(\frac{|\mathcal{X}|}{W} |\mathcal{F}|)$, with
 311 $len = end - ini + 1$, at the worst case with $len \leq \min\{|\mathcal{X}|, |\mathcal{F}|\}$. The expla-
 312 nation is given that we can solve any set operations between a pair of sets
 313 —such as UNION, INTERSECTION and DIFFERENCE— in linear time on the
 314 size $|\mathcal{X}|/W$, where W is the memory word size in the CPU architecture³.
 315 This is achieved with efficient computing techniques that are detailed in Sec-
 316 tion 3.2. Accordingly, we also take $O(|\mathcal{X}|/W)$ time for computing $\mathcal{U} - S$ and

³Note that if $n' = |\mathcal{X}|/W$ is small, then the operation is made in constant time.

317 $\mathcal{C} \cup S$, at lines 13—14. Altogether, the total run time of our algorithm is
 318 $O(\min\{|\mathcal{X}|, |\mathcal{F}|\} \frac{|\mathcal{X}| \cdot |\mathcal{F}|}{W})$ at the worst case, which is clearly polynomial.

319

320 Let us now continue with the detail of the approximation relation of our
 321 solution, which we will proof that it cannot be greater than $|\mathcal{C}^*|(\ln |\mathcal{X}| + 1)$ in
 322 the worst case, where \mathcal{C}^* is an optimal solution for the instance. To perform
 323 the proof, we will use the following inequality that relates the Harmonic
 324 Number $H(n)$ with the natural logarithm [17, 28].

$$H(n) = \frac{1}{1} + \frac{1}{2} + \dots + \frac{1}{n} = \sum_{k=1}^n \frac{1}{k} \leq \ln n + 1 \quad (1)$$

325 Let S_i denote the i th subset selected by our SUCCINCT-SETCOVER Al-
 326 gorithm 2 and k_i the current grade of minimal items at the i th iteration
 327 respect to the inner loop (lines 11—16). Thus, the i th iteration corresponds
 328 when the procedure FINDSUBSET(...) selects the best subset $S_i \in \mathcal{F}$ from
 329 the range of degree k_i , which is the one that contains the largest number of
 330 uncovered minimal elements (i.e., elements of degree k_i). We also establish
 331 that each selected subset S_i contributes with a cost of 1 to \mathcal{C} , cost that is
 332 distributed by all the new minimal items $x \in S_i$, and the sum of all these
 333 costs is also the cardinality of our solution \mathcal{C} . We define the cost of $x \in S_i$
 334 as follow

$$c_x = \begin{cases} \frac{1}{\text{MINIMAL}(S_i - \mathcal{C}_{i-1})} & , \text{ if } x \in S_i \text{ has grade } k_i \text{ and } x \notin \mathcal{C}_{i-1} \\ 0 & , \text{ otherwise} \end{cases} \quad (2)$$

335 where $\mathcal{C}_{i-1} = S_1 \cup S_2 \cup \dots \cup S_{i-1}$ is the partial solution built up to $(i-1)$ th
 336 iteration and $\text{MINIMAL}(S_i - \mathcal{C}_{i-1})$ corresponds to the number of items of grade
 337 k_i added by first time by the algorithm in the i th iteration. Seen another
 338 way, if S_i is selected, then $\text{MINIMAL}(S_i - \mathcal{C}_{i-1})$ is the number of new minimal
 339 elements that S_i contributes to the solution. Thus, we have

$$\sum_{x \in S_i} c_x = 1 \quad (3)$$

340 Now, for any set S belonging to the family \mathcal{F} and by using our definition
 341 of cost the following relation is always true

$$\sum_{x \in S} c_x \leq H(|S|) \quad (4)$$

342 In order to prove (4), consider any set $S \in \mathcal{F}$ and any $i = 1, 2, \dots, |\mathcal{C}|$,
 343 and let $u_i = \text{MINIMAL}(S - (S_1 \cup S_2 \cup \dots \cup S_i))$, as the number of minimal
 344 elements in S that remain uncovered after the algorithm has selected sets
 345 S_1, S_2, \dots, S_i . Let $k' \geq 1$ the minimal grade for any element in the universe
 346 and $u_0 = \text{MINIMAL}(S)$ as the number of elements with grade k' in S . Let
 347 $t \leq |\mathcal{C}|$, the smallest index such that $S \subseteq (S_1 \cup \dots \cup S_{t-1} \cup S_t)$ but $\exists x \in$
 348 $S : x \notin (S_1 \cup \dots \cup S_{t-1})$. Thus, $u_t = 0$ and $u_j \leq u_{j+1}, \forall j = 0, 1, \dots, t-1$.
 349 Observe that $(u_{i-1} - u_i) \geq 0$ minimal elements of S are covered for the first
 350 time by S_i , for $i = 1, 2, \dots, t$. Then we have (for S)

$$\begin{aligned}
 \sum_{x \in S} c_x &= \sum_{i=1}^t (u_{i-1} - u_i) \cdot c_x \\
 &\leq \sum_{i=1}^t (u_{i-1} - u_i) \frac{1}{\text{MINIMAL}(S_i - (S_1 \cup S_2 \cup \dots \cup S_{i-1}))}
 \end{aligned}$$

351 because it is possible that for some elements in S the cost is 0. We can also see
 352 that $\text{MINIMAL}(S_i - (S_1 \cup S_2 \cup \dots \cup S_{i-1})) \geq \text{MINIMAL}(S - (S_1 \cup S_2 \cup \dots \cup S_{i-1}))$,
 353 since the algorithm selects the set S_i as the set with the largest number of
 354 uncovered elements of grade k_i (i.e. $\text{MINIMAL}(S_i) \geq \text{MINIMAL}(S), \forall S \in \mathcal{F}$
 355 not included yet). Then, we now have

$$\begin{aligned}
 \sum_{x \in S} c_x &\leq \sum_{i=1}^t (u_{i-1} - u_i) \frac{1}{\text{MINIMAL}(S - (S_1 \cup S_2 \cup \dots \cup S_{i-1}))} \\
 &\leq \sum_{i=1}^t (u_{i-1} - u_i) \frac{1}{u_{i-1}} \\
 &= \sum_{i=1}^t \sum_{j=u_i+1}^{u_{i-1}} \frac{1}{u_{i-1}}
 \end{aligned}$$

356 giving that $u_{i-1} \geq j$, we have

$$\begin{aligned}
\sum_{x \in S} c_x &\leq \sum_{i=1}^t \sum_{j=u_i+1}^{u_{i-1}} \frac{1}{j} \\
&= \sum_{i=1}^t \left(\sum_{j=1}^{u_{i-1}} \frac{1}{j} - \sum_{j=1}^{u_i} \frac{1}{j} \right) \\
&= \sum_{i=1}^t (H(u_{i-1}) - H(u_i)) \\
&= H(u_0) - H(u_t) \\
&= H(\text{MINIMAL}(S)) - 0 \quad (\text{because } u_0 = \text{MINIMAL}(S)) \\
&\leq H(|S|) \quad (\text{because } \text{MINIMAL}(S) \leq |S|)
\end{aligned}$$

357 We now can prove Eq. (1). Thus, giving our definition of cost, we have

$$|\mathcal{C}| = \sum_{x \in \mathcal{X}} c_x \quad (5)$$

358 As each element $x \in \mathcal{X}$ is at least in one set of the optimal solution \mathcal{C}^* , then
359 we have

$$|\mathcal{C}| = \sum_{x \in \mathcal{X}} c_x \leq \sum_{S \in \mathcal{F}} \sum_{x \in S} c_x \quad (6)$$

360 For (4) we have

$$\begin{aligned}
|\mathcal{C}| &\leq \sum_{S \in \mathcal{F}} |H(S)| \\
&\leq |\mathcal{C}^*| \cdot H(\max\{|S| : S \in \mathcal{F}\})
\end{aligned}$$

361 Finally, by using (1) we conclude that

$$|\mathcal{C}| \leq |\mathcal{C}^*|(\ln |\mathcal{X}| + 1)$$

362 which completes the proof that the approximation ratio of our algorithm is
363 $\ln |\mathcal{X}| + 1$.

3.2. Efficient Computing Techniques in Succinct Space

In order to achieve sublinear time for operations on sets such as UNION, INTERSECTION and DIFFERENCE, we first represent the input datasets $(\mathcal{X}, \mathcal{F})$ as a bit matrix $M_{|\mathcal{F}| \times |\mathcal{X}|}$ of $|\mathcal{F}| \cdot |\mathcal{X}|$ bits. Thus, the i th row $M[i]$ in $M_{|\mathcal{F}| \times |\mathcal{X}|}$ corresponds to the representation of the i th subset S_i , and the j th element of the universe \mathcal{X} is represented by the j th bit in each row of the matrix. We then set $M[i][j] = 1$ if the j th element is present in the subset S_i and with 0 if this element does not belong to the set. We now explain $\forall i, j, 1 \leq i \leq |\mathcal{X}|, 1 \leq j \leq |\mathcal{F}|$, how to compute $S_i \cup S_j$, $S_i \cap S_j$ or $S_i - S_j$ in constant time from the bit representation.

First, consider the singular case when $|\mathcal{X}| \leq W$, where W is the memory word size of the CPU architecture. Let p and q be two words of memory, of W bits each, to represent the $\leq W$ elements of the sets S_i and S_j respectively. Therefore, we can solve any of these set operations on p and q simply by applying the *bitwise* operators AND, OR and NEGATION. The table 1 shows how to compute these operations between sets by using bitwise operator.

set operation	bitwise operation
$S_i \cup S_j$	$p \text{ OR } q$
$S_i \cap S_j$	$p \text{ AND } q$
$S_i - S_j$	$p \text{ AND } \neg q$

Table 1: The equivalent bitwise operations for set operations. The sets S_i and S_j are represented by the bit sequences p and q respectively, whose lengths are at most $|W|$ bits (i.e., $|\mathcal{X}| \leq W$ elements). The symbol \neg corresponds to the negation operator (i.e., $\neg 1 = 0$ and $\neg 0 = 1$). Therefore, to compute any set operation can be done in $O(1)$ by using bitwise operator on p and q .

For universes of any length, $|\mathcal{X}| > W$, we can process the $|\mathcal{X}|$ bits by implement a loop routine in $O(|\mathcal{X}|/W)$ time; where p and q are bit arrays of $|\mathcal{X}|/W$ words of memory.

385 4. Experimental Results

386 In order to evaluate the empiric performance of the algorithms, we pro-
387 vide three implementations in C++⁴. The first, called **Exhaustive-SC**, finds
388 the optimal set-covering by means of exhaustive searches. This is used as a
389 baseline for the smaller datasets described in the experiment 2. The second,
390 referred to as **Greedy-SC**, is the implementation of the classical approximate
391 greedy set cover solution given by the Algorithm 1. The last implementa-
392 tion corresponds to our **SUCCINCT-SETCOVER** algorithm which is referred
393 to as **Succinct-SC**. We test the implementations by using different instances
394 $(\mathcal{X}, \mathcal{F})$ of different lengths and different density for each set. All runtimes
395 obtained are shown in logarithmic scale.

396 4.1. Setup

397 The experiments mentioned were performed in the Patagón Supercom-
398 puter [30], using a compute node with $2 \times$ AMD EPYC 7742 CPU (2.6Ghz,
399 64-cores, 256MB L3 Cache), equipped with 1TB RAM DDR4-3200Hz. Only
400 a single thread of execution was used. The OS was Linux (Ubuntu 20.04,
401 64bit) running kernel 5.8.0. All programs were compiled using **g++** version
402 9.3.0. The runtime corresponds to real (wallclock) time from the **Chrono**
403 library, using the **high.resolution** clock method.

404 4.2. Experiments with Random-Generated Datasets

405 In the first cycle of experiments, Figure 2, we compute the exact solution
406 for the MSCP and the approximate responses to each instance. Since the
407 time complexity of the optimal solution is superpolynomial, we only run
408 experiments for instances small enough to get the answer in a reasonable
409 amount of time. The experiment involves 4 cases with different sizes of
410 $|\mathcal{F}|$, and each one with 5 instances with an expected universe of $|\mathcal{X}| = 70$,
411 randomly generated using a Gaussian distribution, whose arithmetic mean
412 was obtained from a uniform distribution over the range $[0..|\mathcal{X}|]$ and the
413 standard deviation is $0.125|\mathcal{X}|$.

414 The dark gray line gave by the **Exhaustive-SC** corresponds to the min-
415 imum cardinality to the problem and therefore, the lower bound for any
416 approximation to the problem.

⁴The repository with the code is available at https://github.com/kokee07/succinct_mscp

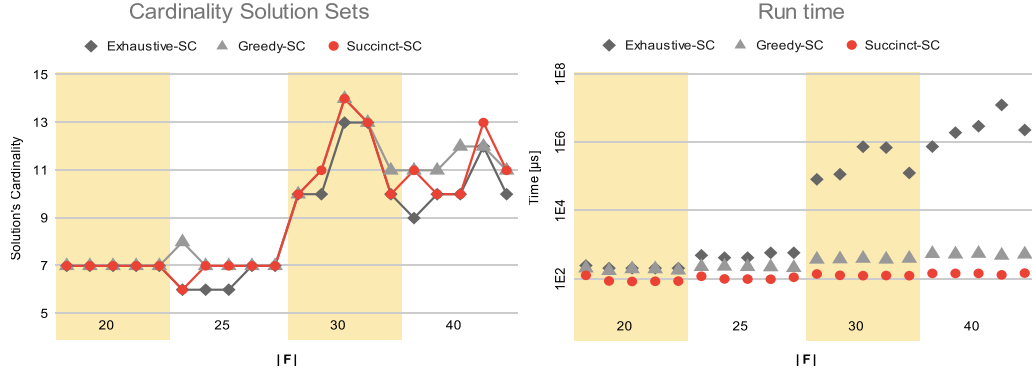


Figure 2: Cardinality (left) and running time (right) of our **Succinct-SC** approximation compared to the exact solution, **Exhaustive-SC**, and the classical greedy heuristic, **Greedy-SC**, for low size random generated datasets.

417 We can appreciate how, on average, the **Succinct-SC** is faster (right) and
 418 smaller (left) than the **Greedy-SC** approximation.
 419 Given the exponential growth in time required for the **Exhaustive-SC** to be
 420 computed on large datasets, the following experiments will only consider the
 421 approximate solutions. The second cycle of experiments, Figures 3 to 5, in-
 422 volves 9 random-generated cases, using the same data generator mentioned
 423 before, but increasing the value of $|\mathcal{X}|$ and $|\mathcal{F}|$ progressively. Each figure
 424 includes 3 cases with 5 instances each, grouped by the size of $|\mathcal{F}|$.

425

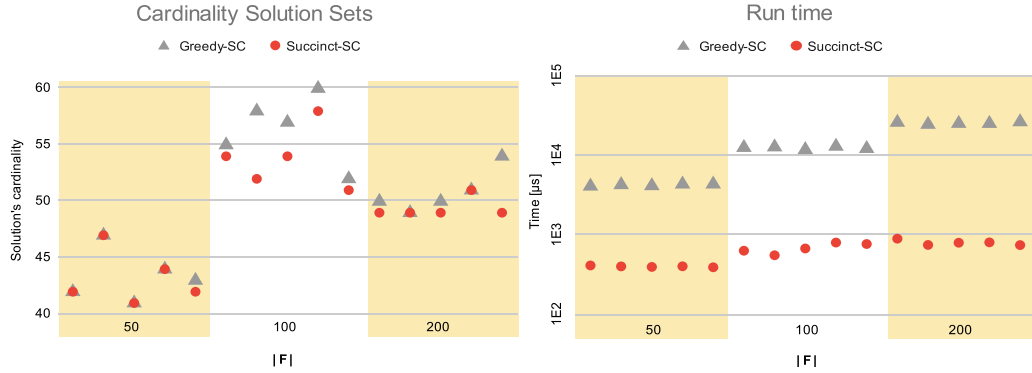


Figure 3: Cardinality and running time of our **Succinct-SET COVER** algorithm compared to the **GREEDY SET COVER** for random-generated datasets with values: $|\mathcal{X}| \in [200..300]$ and $|\mathcal{F}| \in [50, 100, 200]$.

Figure 3 (left) shows an irregular difference between the quality of the solutions of both implementations, being **Succinct-SC** at least as good as **Greedy-SC**. Regarding time (right), as $|\mathcal{F}|$ increases, the distance between the execution time of each algorithm keeps increasing, at least by an order of magnitude.

Figure 4 shows the results over a medium size dataset. Achieving a more regular distance between the solutions quality, being our algorithm slightly better than the original Greedy approach. At the right side, each case of $|\mathcal{F}|$ size has an increase in the order of magnitude between the both solutions, evidencing the higher complexity of operations required for **Greedy-SC**.

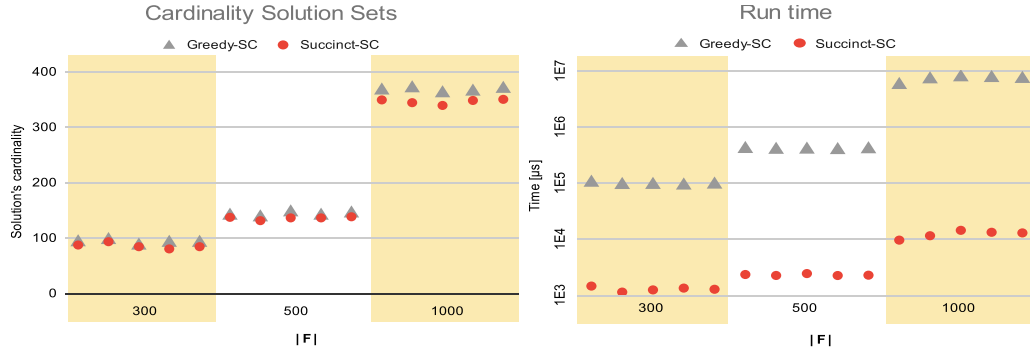


Figure 4: Cardinality and running time of our **SUCCINCT-SET COVER** algorithm compared to the **GREEDY SET COVER** for random-generated datasets with values: $|\mathcal{X}_e| \in [500..2.000]$ and $|\mathcal{F}| \in [300, 500, 1.000]$.

In the Figure 5, we generate higher order of $|\mathcal{F}|$. The left side shows the cardinality of the solution sets in thousands of subsets, in which our **Succinct-SC** algorithm computes a lower, and hence better, solution for each instance than the **Greedy-SC**. At the right side, the distance in order of magnitude keeps the same behaviour as the previous experiments.

The next cycle of experiments measures the advantage of our algorithm for the case where there are elements that are included by a single subset of \mathcal{F} , or by very few sets of \mathcal{F} . In this scenario, our algorithm takes advantage by performing, in the first iteration cycles, a much more efficient pruning of the search space than the traditional greedy approach. To emulate that scenario, we implement a new generator to secure the existence of elements contained only by one subset of $|\mathcal{F}|$. To achieve this, the generator takes

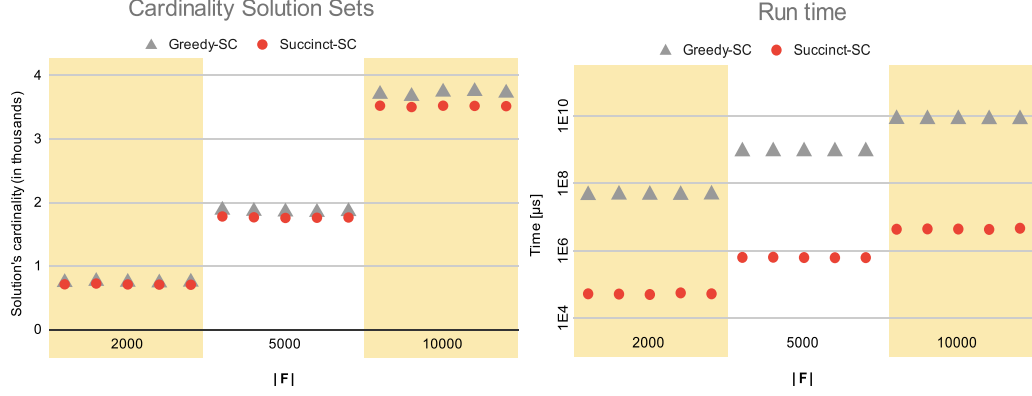


Figure 5: Cardinality and running time of our SUCCINCT-SET COVER algorithm compared to the GREEDY SET COVER for random-generated datasets with values: $|\mathcal{X}| \in [4.000..20.000]$ and $|\mathcal{F}| \in [2.000, 5.000, 10.000]$.

450 as input the expected universe size $|\mathcal{X}|$, the total number of subsets $|\mathcal{F}|$, a
 451 number \mathcal{T} of samples to generate and a parameter p . Thus, we ensure that
 452 the first $p * |\mathcal{X}|$ elements of the universe will be randomly assigned to one set
 453 each; and then we generate all the rest of the \mathcal{T} samples by using a uniform
 454 distribution over the range $[p * |\mathcal{X}| \dots |\mathcal{X}|]$.

455

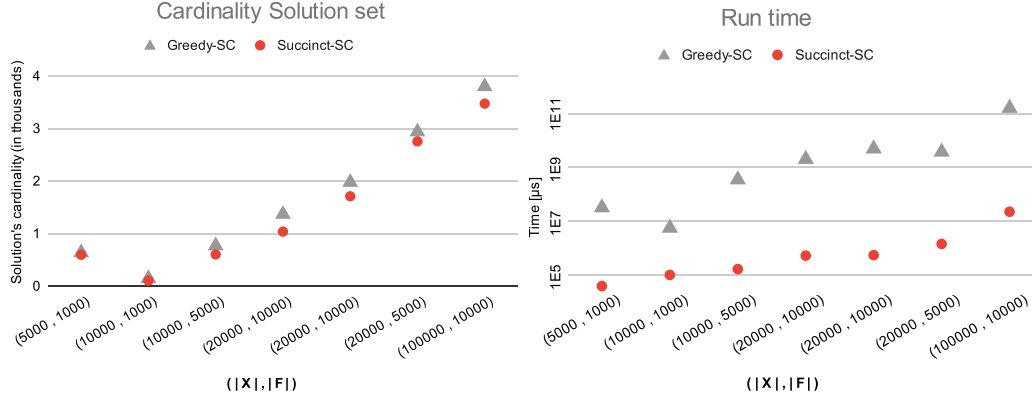


Figure 6: Cardinality and running time of our SUCCINCT-SET COVER algorithm compared to the GREEDY SET COVER for optimal random-generated datasets.

456 Figure 6 shows the results obtained from 7 different instances of optimal
 457 datasets, with the pair values $(|\mathcal{X}|, |\mathcal{F}|)$ of each case indicated in the X-

axis. The results show that **Succinct-SC** obtains a better solution for any instance within this scenario, with a greater difference than the previous cases. Furthermore, as shown on the right side of Figure 6, in these cases a difference of up to 4 orders of magnitude is reached between the times of both algorithms.

4.3. Experiments with Natural/Real Datasets

The following experiments show the results of Dataset benchmarks from OR-Library[9], including the **SET PARTITIONING PROBLEMS** datasets from K. L. Hoffman and D. Levine[8], used in genetic algorithms, in Figure 7 and **SET COVER PROBLEM** datasets from a Railway System, contributed by Paolo Nobili[7] in Figure 8.

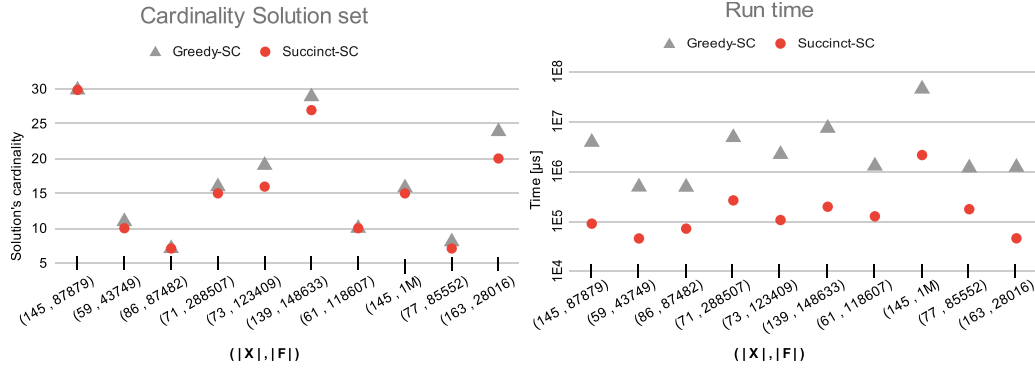


Figure 7: Cardinality and running time of our **SUCCINCT-SET COVER** algorithm compared to the **GREEDY SET COVER** for Genetic Datasets with a low \mathcal{F} cardinality.

The decision of using the genetic algorithm Dataset as a Natural case is supported due to the low universe cardinality with a high number of combinations from \mathcal{F} presented. For our purpose, the constraints associated to **SET PARTITIONING PROBLEMS** in the dataset are ignored, such as that any element must be included in one and only one set in the solution. In Figure 7 we can see, at the left, a more erratic behaviour from each case due to their natural origin. At the right, our **SUCCINCT-SET COVER** is at least 1 order of magnitude faster than the greedy approach.

As the last experiment, we test the **SET COVER PROBLEM** solution for a larger dataset from Italian railways [7], where a set cover solution represents a selection of routes that cover all the stations. The Figure 8 shows the results from 7 cases, indicating their respective $|\mathcal{X}|$ and $|\mathcal{F}|$ at the X-axis.

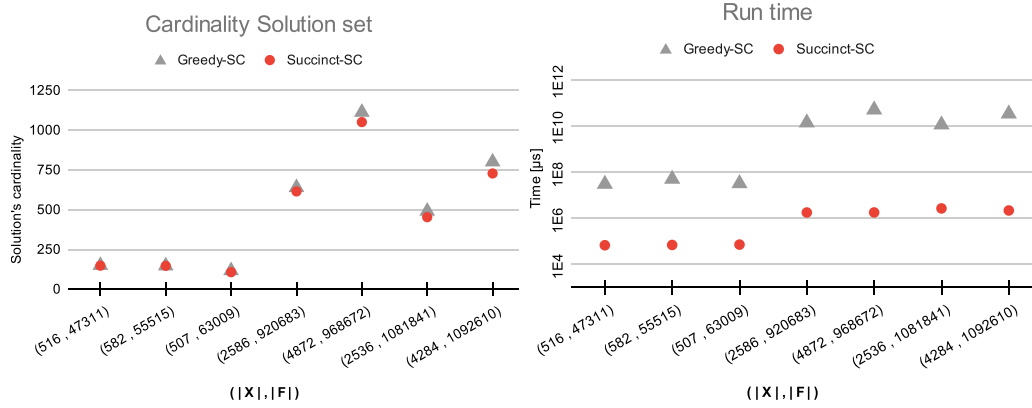


Figure 8: Cardinality and running time of our **SUCCINCT-SET COVER** algorithm compared to the **GREEDY SET COVER** for Railways Datasets.

482 The results shown in Figure 8 are consistent with the all the previous
 483 experiments, with our **SUCCINCT-SET COVER** algorithm achieving a solution
 484 at least as good as **GREEDY SET COVER** for those cases with a low universe,
 485 and a better solution for the other cases with a high universe of elements. At
 486 the right side, the difference in order of magnitude between both algorithm
 487 varies from 2 to 3 orders in each case. This is highly important since our
 488 technique could find a better solution in minutes for any real case, in contrast
 489 with the hours or days of processing attributed to the Greedy algorithm.

490 5. Conclusions and Future Work

491 We have presented a polynomial time approximate solution to solve the
 492 MSCP, offering an innovative algorithm that includes data compression tech-
 493 niques to speed up the process. The proposed **Succinct-SC** algorithm pre-
 494 processes the input dataset to create a succinct representation of the sets of
 495 the family \mathcal{F} ; which also allows us to calculate operations between sets in a
 496 much more efficient way than with the traditional approach. Thanks to this,
 497 we improve the execution time of the **Greedy-SC** algorithm significantly, both
 498 in theory and in practice. This theoretical asymptotic bound is an important
 499 part of the contributions of this work. The algorithm **Succinct-SC** also offers
 500 a theoretically optimal ratio approximation under the assumption that there
 501 is no polynomial time algorithm that solves any of the problems belonging
 502 to the NP-Hard class. We have implemented our algorithm to measure the
 503 computation time and the quality of the approximation, demonstrating that

504 Succinct-SC improves the query time by orders of magnitude with respect to
 505 to the traditional greedy approximation, Greedy-SC algorithm, while achieving
 506 a better approximation rate. In summary, the results obtained with our
 507 algorithm are very favorable for the datasets tested, which included randomly
 508 generated ones as well as real ones from the state of the art.

509

510 One aspect to highlight is in relation to the efficiency achieved thanks to
 511 the succinct representation of the data, which allows computing operations
 512 between sets at high speed, always in a succinct space. We establish that
 513 the main focus when using this technique is to achieve greater speed and not
 514 necessarily compression. Since we are working on an approximation, which
 515 first has to guarantee a good quality response—we achieve this with a theo-
 516 retical proof of the bound of the algorithm—but at high speed—we achieve
 517 it by providing an implementation of an algorithm polynomial that also, in
 518 practice, solves bitwise set operations.

519

520 The succinct/compressed representation of the data is a relevant sub-
 521 ject for future work, where the goal will be to design an optimal algorithm
 522 that works in a succinct space, also incorporating CPU/GPU parallelism.
 523 In terms of compression and based on the tested scenarios, we expect that
 524 our implementation can reach volumes that are difficult for other implemen-
 525 tations to reach; standing out even more in universes with a large number
 526 of subsets (when \mathcal{F} is very large), where each subset has a large number of
 527 elements on average.

528 6. Acknowledgment

529 This research was supported by the Fondecyt grant #11221029 and by
 530 the Patagón supercomputer of Universidad Austral de Chile (FONDEQUIP
 531 EQM180042).

532 References

- 533 [1] Balas, E. and Carrera, M. C. (1996). A dynamic subgradient-based
 534 branch-and-bound procedure for set covering. *Oper. Res.*, 44(6):875–890.
- 535 [2] Balas, E. and Ho, A. (1980). *Set covering algorithms using cutting planes,
 536 heuristics, and subgradient optimization: A computational study*, pages 37–
 537 60. Springer Berlin Heidelberg, Berlin, Heidelberg.

- 538 [3] Balas, E. and Ho, A. (2009). *Set covering algorithms using cutting*
539 *planes, heuristics, and subgradient optimization: A computational study*,
540 volume 12, pages 37–60.
- 541 [4] Barnhart, C., Cohn, A. M., Johnson, E. L., Klabjan, D., Nemhauser,
542 G. L., and Vance, P. H. (2003). *Airline Crew Scheduling*, pages 517–560.
543 Springer US, Boston, MA.
- 544 [5] Beasley, J. (1987). An algorithm for set covering problem. *European*
545 *Journal of Operational Research*, 31(1):85–93.
- 546 [6] Beasley, J. and Jørnsten, K. (1992). Enhancing an algorithm for set
547 covering problems. *European Journal of Operational Research*, 58(2):293–
548 300. Practical Combinatorial Optimization.
- 549 [7] Beasley, J. E. Or-library, set cover problems dataset. [dataset].
- 550 [8] Beasley, J. E. Or-library, set partitioning problems dataset. [dataset].
- 551 [9] Beasley, J. E. (1990). Or-library: Distributing test problems by electronic
552 mail. *Journal of the Operational Research Society*, 41(11):1069–1072.
- 553 [10] Bilal, N., Galinier, P., and Guibault, F. (2013). A new formulation of
554 the set covering problem for metaheuristic approaches. *ISRN Operations*
555 *Research*, 2013:10 pages.
- 556 [11] Burhan, H. M., Attea, B. A., Abbood, A. D., Abbas, M. N., and Al-
557 Ani, M. (2021). Evolutionary multi-objective set cover problem for task
558 allocation in the internet of things. *Applied Soft Computing*, 102:107097.
- 559 [12] Caprara, A., Fischetti, M., and Toth, P. (1999). A heuristic method for
560 the set covering problem. *Operations Research*, 47(5):730–743.
- 561 [13] Chan, T. M., Grant, E., Könemann, J., and Sharpe, M. (2012).
562 Weighted capacitated, priority, and geometric set cover via improved quasi-
563 uniform sampling. In *Proceedings of the Twenty-Third Annual ACM-SIAM*
564 *Symposium on Discrete Algorithms*, SODA ’12, page 1576–1585, USA. So-
565 ciety for Industrial and Applied Mathematics.
- 566 [14] Chen, J., Mi, J., and Lin, Y. (2020). A graph approach for fuzzy-rough
567 feature selection. *Fuzzy Sets and Systems*, 391:96–116. Computer Science.

- 568 [15] Chvatal, V. (1979). A greedy heuristic for the set-covering problem.
569 *Math. Oper. Res.*, 4(3):233–235.
- 570 [16] Contardo, C. and Hertz, A. (2021). An exact algorithm for a class of
571 geometric set-cover problems. *Discrete Applied Mathematics*, 300:25–35.
- 572 [17] Cormen, T. H., Leiserson, C. E., Rivest, R. L., and Stein, C. (2009).
573 *Introduction to Algorithms, 3rd Edition*. MIT Press.
- 574 [18] Feige, U. (1998). A threshold of $\ln n$ for approximating set cover. *J.*
575 *ACM*, 45(4):634–652.
- 576 [19] Fisher, M. L. and Kedia, P. (1990). Optimal solution of set
577 covering/partitioning problems using dual heuristics. *Manage. Sci.*,
578 36(6):674–688.
- 579 [20] Gomes, F. C., Meneses, C. N., Pardalos, P. M., and Viana, G. V. R.
580 (2006). Experimental analysis of approximation algorithms for the ver-
581 tex cover and set covering problems. *Computers & Operations Research*,
582 33(12):3520–3534. Part Special Issue: Recent Algorithmic Advances for
583 Arc Routing Problems.
- 584 [21] Johnson, D. S. (1973). Approximation algorithms for combinatorial
585 problems. In *Proceedings of the Fifth Annual ACM Symposium on Theory*
586 *of Computing*, STOC '73, page 38–49, New York, NY, USA. Association
587 for Computing Machinery.
- 588 [22] Karp, R. M. (1972). *Reducibility among Combinatorial Problems*, pages
589 85–103. Springer US, Boston, MA.
- 590 [23] Knuth, D. E. (1973). *The Art of Computer Programming. Vol.3: Sorting*
591 *and Searching*.
- 592 [24] Koufogiannakis, C. and Young, N. E. (2009). Distributed and parallel
593 algorithms for weighted vertex cover and other covering problems. In *Pro-*
594 *ceedings of the 28th ACM Symposium on Principles of Distributed Com-*
595 *puting*, PODC '09, page 171–179, New York, NY, USA. Association for
596 Computing Machinery.
- 597 [25] Lan, G., DePuy, G. W., and Whitehouse, G. E. (2007). An effective
598 and simple heuristic for the set covering problem. *European Journal of*
599 *Operational Research*, 176(3):1387–1403.

- [26] Lim, C. L., Moffat, A., and Wirth, A. (2014). Lazy and eager approaches for the set cover problem. In *Proceedings of the Thirty-Seventh Australasian Computer Science Conference - Volume 147*, ACSC '14, page 19–27, AUS. Australian Computer Society, Inc.
- [27] Lund, C. and Yannakakis, M. (1994). On the hardness of approximating minimization problems. *J. ACM*, 41(5):960–981.
- [28] Olaikhan, A. (2021). *An Introduction to the Harmonic Series and Logarithmic Integrals: For High School Students Up to Researchers*. Amazon Digital Services LLC - KDP Print US.
- [29] Owais, M., Osman, M. K., and Moussa, G. (2016). Multi-objective transit route network design as set covering problem. *IEEE Transactions on Intelligent Transportation Systems*, 17(3):670–679.
- [30] Patagón Supercomputer (2021). <https://patagon.uach.cl>. Accessed January 21, 2023.
- [31] Peleg, D., Schechtman, G., and Wool, A. (1997). Randomized approximation of bounded multicovering problems. *Algorithmica*, 18(1):44–66.
- [32] Peuzin-Jubert, M., Polette, A., Nozais, D., Mari, J.-L., and Pernot, J.-P. (2021). Survey on the view planning problem for reverse engineering and automated control applications. *Computer-Aided Design*, 141:103094.
- [33] Ribeiro, C. C., Minoux, M., and Penna, M. C. (1989). An optimal column-generation-with-ranking algorithm for very large scale set partitioning problems in traffic assignment. *European Journal of Operational Research*, 41(2):232–239.
- [34] Spensieri, D., Åblad, E., Bohlin, R., Carlson, J. S., and Söderberg, R. (2021). Modeling and optimization of implementation aspects in industrial robot coordination. *Robotics and Computer-Integrated Manufacturing*, 69:102097.
- [35] Wang, X. and Ju, S. (2009). A set-covering-based approach for overlapping resource selection in distributed information retrieval. In *2009 WRI World Congress on Computer Science and Information Engineering*, volume 4, pages 272–276.