

Set Cover

Lecturer: Seth Gilbert

August 25, 2015

Abstract

Here we consider the problem of set cover. We begin with an example, and then discuss the classical greedy solution. Finally, we analyze the greedy algorithm and show that it is an $O(\log n)$ -approximation algorithm. In the analysis, notice how we begin by giving a lower bound on OPT, showing that OPT can only do so well. We then relate this to the performance of the Greedy Algorithm to get the approximation ratio.

1 Set Cover

The problem of *Set Cover* is a combinatorial optimization problem that frequently shows up in real-world scenarios, typically when you have collections of needs (e.g., tasks, responsibilities, or capabilities) and collections of resources (e.g., employees or machines) and you need to find a minimal set of resources to satisfy your needs. In fact, Set Cover generalizes Vertex Cover: in vertex cover, each vertex (i.e., set) covers the adjacent edges (i.e., elements); in set cover, each set can cover an arbitrary set of elements.

Unfortunately, Set Cover is NP-hard (i.e., NP-complete as a decision problem). As an exercise, prove that Set Cover is NP-hard by reducing another NP-complete problem to Set Cover.

In fact, we know that Set Cover is hard to approximate. There are no polynomial-time α -approximation algorithms for any constant, assuming $P \neq NP$. In fact, if there were a polynomial-time $(1 - \epsilon) \ln n$ -approximation algorithms, for any $\epsilon > 0$, then we could solve every problem in NP in $O(n^{\log \log n})$ time. Most computer scientists believe that this is unlikely!

1.1 Problem Definition

We first define the problem, and then give some examples that show how set cover might appear in the real world.

Definition 1 Let $X = \{x_1, x_2, \dots, x_n\}$ be a set of n elements. Let S_1, S_2, \dots, S_m be subsets of X , i.e., each $S_j \subseteq X$. Assume that every item in X appears in some set, i.e., $\bigcup_j S_j = X$. A **set cover** of X with S is a set $I \subseteq \{1, \dots, m\}$ such that $\bigcup_{j \in I} S_j = X$. A **minimum set cover** is a set cover I of minimum size.

That is, a minimum set cover is the smallest set of sets $\{S_{i_1}, S_{i_2}, \dots, S_{i_k}\}$ that covers X .

Example 1. Assume you have a set of software developers: Alice, Bob, Collin, and Dave. Each programmer knows at least one programming language. Alice knows C and C++. Bob knows C++ and Java. Collin knows C++, Ruby, and Python. Dave knows C and Java. Your job is to hire a team of programmers. You are given two requirements: (i) there has to be at least one person on the team who knows each language (i.e., C, C++, Java, Python, and Ruby), and (ii) your team should be as small as possible.

This is precisely a set cover problem. The base elements X are the 5 different programming languages. Each programmer represents a set. Your job is to find the minimum number of programmers (i.e., the minimum number of sets) such that every language is covered.

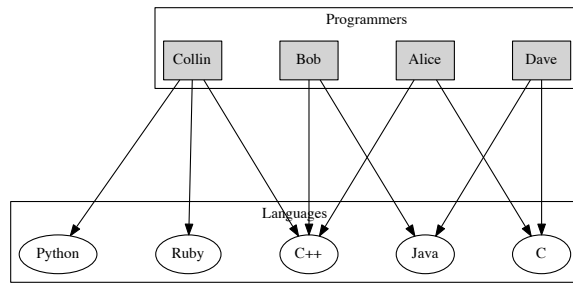


Figure 1: The problem of hiring software developers, represented as a bipartite graph. The goal is to choose a minimum sized set of programmers to ensure that every language is known by at least one of your programmers.

See Figure 1 for an example of this problem, depicted here as a bipartite graph. You will notice that any set cover problem can be represented as a bipartite graph, with the sets represented on one side and the base elements represented on the other side.

In this case, Alice, Bob, and Collin form a set cover of size 3. However, Collin and Dave form a set cover of size 2, which is optimal.

Example 2. Any vertex cover problem can be represented as a set cover problem. Assume you are given a graph $G = (V, E)$ and you want to find a vertex cover. In this case, you can define $X = E$, i.e., the elements to be covered are the edges. Each vertex represents a set. We define the set $S_u = \{(u, v) : (u, v) \in E\}$, i.e., S_u is the set of edges adjacent to u . The problem of vertex cover is now to find a set of sets $S_{i_1}, S_{i_2}, \dots, S_{i_k}$ such that every edge is covered.

1.2 Greedy Algorithm

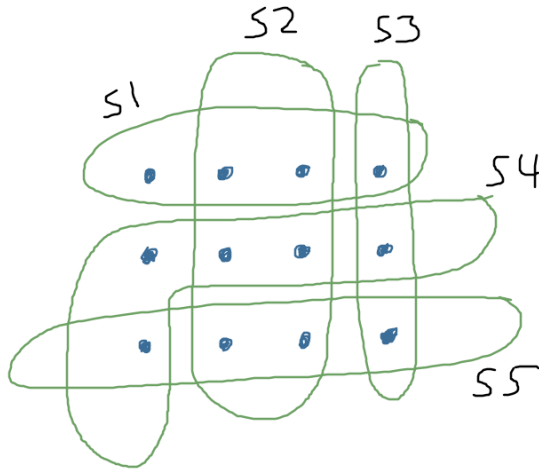
There is a simple greedy algorithm for solving set cover:

```

/* This algorithm adds sets greedily, one at a time, until everything is
   covered. At each step, the algorithm chooses the next set that will
   cover the most uncovered elements. */
1 Algorithm: GreedySetCover( $X, S_1, S_2, \dots, S_m$ )
2 Procedure:
3    $I \leftarrow \emptyset$ 
   /* Repeat until every element in  $X$  is covered: */
4   while  $X \neq \emptyset$  do
5     Let  $d(j) = |S_j \cap X|$  // This is the number of uncovered elements in  $S_j$ 
6     Let  $j = \operatorname{argmax}_{i \in \{1, \dots, m\}} d(i)$ 
7      $I \leftarrow I \cup \{j\}$ 
8      $X \leftarrow X \setminus S_j$  // Remove elements in  $S_j$  from  $X$ .
9   return  $I$ 

```

The algorithm proceeds greedily, adding one set at a time to the set cover until every element in X is covered by at least one set. In the first step, we add the set that covers the most elements. At each ensuing step, the algorithm chooses the set that covers the most elements that remain uncovered.



| Set | Step 1 | Step 2 | Step 3 | Step 4 | Step 5 |
|-------|--------|--------|--------|--------|--------|
| S_1 | 4 | 2 | 1 | 1 | 0 |
| S_2 | 6 | 0 | 0 | 0 | 0 |
| S_3 | 3 | 3 | 0 | 0 | 0 |
| S_4 | 5 | 3 | 2 | 0 | 0 |
| S_5 | 4 | 2 | 1 | 0 | 0 |

Figure 2: On the left is an example of set cover consisting of twelve elements and five sets. On the right is a depiction of what happens when you execute the GreedySetCover algorithm on this example. Each column represents the number of new elements covered by each set at the beginning of the step.

Let's look at an example—see Figure 2. Here we have 12 elements (represented by the dots) and 5 sets: S_1, S_2, S_3, S_4, S_5 . In the first iteration, we notice that set S_2 covers the most elements, i.e., 6 elements, and hence it is added to the set cover. In the second iteration, set S_3 and S_4 both cover 3 new elements, and so we add set S_3 to the set cover. In the third iteration, set S_4 covers 2 new elements, and so we add it to the set cover. Finally, in the fourth step, set S_1 covers one new element and so we add it to the set cover. Thus, we end up with a set cover consisting of S_1, S_2, S_3, S_4 . Notice, though, that the optimal set cover consists of only three elements: S_1, S_4, S_5 .

1.3 Analysis

Our goal is to show that the greedy set cover algorithm is an $O(\log n)$ approximation of optimal. As is typical, in order to show that it is a good approximation of optimal, we need some way to bound the optimal solution. Throughout this section, we will let OPT refer to the optimal set cover.

Intuition. To get some intuition, let's consider Figure 2 and see what we can say about the optimal solution. Notice that there are 12 elements that need to be covered, and none of the sets cover more than 6 elements. Clearly, then, any solution to the set cover problem requires at least $12/6 = 2$ sets. Now consider the situation after the first iteration, i.e., after adding set S_3 to the set cover. At this point, there are 6 elements that remain to be covered, and none of the sets cover more than 3 elements. Any solution to the set cover problem requires at least $6/3 = 2$ sets.

In general, if at some point during the greedy algorithm, there are only k elements that remain uncovered and none of the sets covers more than t elements, then we can conclude that $OPT \geq k/t$. We will apply this intuition to show that the greedy algorithm is a good approximation of OPT .

Definitions. Assume we run the greedy set cover algorithm on elements X and sets S_1, S_2, \dots, S_m . When we run the algorithm, let us label the elements in the order that they are covered:

$$\underbrace{x_1, x_2, x_3, x_4}_{S_5}, \underbrace{x_5, x_6, x_7}_{S_3}, \underbrace{x_8, x_9}_{S_1}, \underbrace{x_{10}}_{S_8}$$

That is, x_1 is the first element covered, x_2 is the second element covered, etc. Under each element, I have indicated the the first set that covered it. In this example, notice that the first set chosen (S_5) covers 4 new elements, the second set chosen (S_3) covers 3 new elements, etc. Each successive set covers at most the same number of elements as the previous one, because the algorithm is greedy: for example, if S_1 here had covered more new elements than S_3 , then it would have been selected before S_3 .

For each element x_j , let c_j be the number of elements covered at the same time. In the example above, this would yield:

$$c_1 = 4, c_2 = 4, c_3 = 4, c_4 = 4, c_5 = 3, c_6 = 3, c_7 = 3, c_8 = 2, c_9 = 2, c_{10} = 1$$

We define $cost(x_j) = 1/c_j$. In this way, the cost of covering all the new elements for some set is exactly 1. In this example, the cost of covering x_1, x_2, x_3, x_4 is 1, the cost of covering x_5, x_6, x_7 is 1, etc. In general, if I is the set cover constructed by the greedy algorithm, then:

$$|I| = \sum_{j=1}^n cost(x_j).$$

Notice that $c_1 \geq c_2 \geq c_3 \geq \dots$, because the algorithm is greedy.

Key step. Let's consider the situation after elements x_1, x_2, \dots, x_{j-1} have been covered already, and the elements x_j, x_{j+1}, \dots remain to be covered. Let OPT be the optimal solution for covering all n elements.

What is the best that OPT can do to cover the element x_j, \dots, x_n ? How many sets does OPT need to cover these remaining elements?

Notice that there remain $n - j + 1$ uncovered elements. However, no set covers more than $c(j)$ of the remaining elements. In particular, all the sets already selected by the greedy algorithm cover *zero* of the remaining elements. Of the sets not yet chosen by the greedy algorithm, the one that covers the most remaining elements covers $c(j)$ of those elements: otherwise, the greedy algorithm would have chosen a different set.

Therefore, OPT needs at least $(n - j + 1)/c(j)$ sets to cover the remaining $(n - j + 1)$ elements. We thus conclude that:

$$OPT \geq \frac{n - j + 1}{c(j)} \geq (n - j + 1)cost(j)$$

Or to put it differently:

$$cost(j) \leq \frac{OPT}{(n - j + 1)}$$

We can now show that the greedy algorithm provides a good approximation:

$$\begin{aligned} |I| &= \sum_{j=1}^n cost(x_j) \\ &\leq \sum_{j=1}^n \frac{OPT}{(n - j + 1)} \\ &\leq OPT \sum_{i=1}^n \frac{1}{i} \\ &\leq OPT(\ln n + O(1)) \end{aligned}$$

(Notice that the third inequality is simply a change of variable where $i = (n - j + 1)$, and the fourth inequality is because the harmonic series $1 + 1/2 + 1/3 + 1/4 + \dots + 1/n$ can be bounded by $\ln n + O(1)$.)

We have therefore shown that the set cover constructed is at most $O(\log n)$ times optimal, i.e., the greedy algorithm is an $O(\log n)$ -approximation algorithm:

Theorem 2 *The algorithm GreedySetCover is a $O(\log n)$ -approximation algorithm for Set Cover.*

There are two main points to note about this proof. First, the key idea was to (repeatedly) bound OPT, so that we could relate the performance of the greedy algorithm to the performance of OPT. Second, the proof crucially depends on the fact that the algorithm is *greedy*. (Always try to understand how the structure of the algorithm, in this case the greedy nature of the algorithm, is used in the proof.) The fact that the algorithm is greedy leads directly to the bound on OPT by limiting the maximum number of new elements that any set could cover.

1.4 Questions

There are many further questions to explore related to set cover. (See the problem set for more.) Here are some natural questions.

1. What is the relationship between Vertex Cover and Set Cover? Show that Algorithm 3 (from Lecture 3) for Vertex Cover is a $O(\log n)$ -approximation algorithm.
2. We have looked here at the *unweighted* versions of set cover: every set has the same cost. Assume instead that each set S_j has a weight $w(S_j)$. The goal is now to minimize the total weight of all the selected sets. Give an approximation algorithm for this weighted version.
3. Give an integer linear program (ILP) for solving the set cover problem.
4. Show how to round the ILP in order to achieve a $O(\log n)$ -approximation of optimal. You may use a randomized algorithm and show that the approximation ratio is achieved with probability at least $1 - 1/n$. (Hint: if you have a variable y_j associated with set S_j , think about what happens when you add set S_j with probability y_j .)
5. Assume that each element x_j appears in at most f sets. Show how to round the ILP in order to achieve an f -approximation. (In this case, use a deterministic rounding scheme.)