

# Greedy Set-Cover Algorithms

(1974-1979, Chvátal, Johnson, Lovász, Stein)

Neal E. Young, University of California, Riverside  
www.cs.ucr.edu/~neal  
entry editor:

**INDEX TERMS:** dominating set, greedy algorithm, hitting set, set cover, minimizing a linear function subject to a submodular constraint

**SYNONYMS:** dominating set, greedy algorithm, hitting set, set cover, minimizing a linear function subject to a submodular constraint

## 1 PROBLEM DEFINITION

Given a collection  $\mathcal{S}$  of sets over a universe  $U$ , a *set cover*  $C \subseteq \mathcal{S}$  is a subcollection of the sets whose union is  $U$ . The *set-cover problem* is, given  $\mathcal{S}$ , to find a minimum-cardinality set cover. In the *weighted set-cover problem*, for each set  $s \in \mathcal{S}$  a weight  $w_s \geq 0$  is also specified, and the goal is to find a set cover  $C$  of minimum total weight  $\sum_{s \in C} w_s$ .

Weighted set cover is a special case of *minimizing a linear function subject to a submodular constraint*, defined as follows. Given a collection  $\mathcal{S}$  of objects, for each object  $s$  a non-negative weight  $w_s$ , and a non-decreasing submodular function  $f : 2^{\mathcal{S}} \rightarrow \mathbb{R}$ , the goal is to find a subcollection  $C \subseteq \mathcal{S}$  such that  $f(C) = f(\mathcal{S})$  minimizing  $\sum_{s \in C} w_s$ . (Taking  $f(C) = |\cup_{s \in C} s|$  gives weighted set cover.)

## 2 KEY RESULTS

The *greedy algorithm* for weighted set cover builds a cover by repeatedly choosing a set  $s$  that minimize the weight  $w_s$  divided by number of elements in  $s$  not yet covered by chosen sets. It stops and returns the chosen sets when they form a cover:

greedy-set-cover( $\mathcal{S}, w$ )

1. Initialize  $C \leftarrow \emptyset$ . Define  $f(C) \doteq |\cup_{s \in C} s|$ .
2. Repeat until  $f(C) = f(\mathcal{S})$ :
3.     Choose  $s \in \mathcal{S}$  minimizing the price per element  $w_s/[f(C \cup \{s\}) - f(C)]$ .
4.     Let  $C \leftarrow C \cup \{s\}$ .
5. Return  $C$ .

Let  $H_k$  denote  $\sum_{i=1}^k 1/i \approx \ln k$ , where  $k$  is the largest set size.

**Theorem 1.** *The greedy algorithm returns a set cover of weight at most  $H_k$  times the minimum weight of any cover.*

*Proof.* When the greedy algorithm chooses a set  $s$ , imagine that it charges the price per element for that iteration to each element newly covered by  $s$ . Then the total weight of the sets chosen by the algorithm equals the total amount charged, and each element is charged once.

Consider any set  $s = \{x_k, x_{k-1}, \dots, x_1\}$  in the optimal set cover  $C^*$ . Without loss of generality, suppose that the greedy algorithm covers the elements of  $s$  in the order given:  $x_k, x_{k-1}, \dots, x_1$ . At the start of the iteration in which the algorithm covers element  $x_i$  of  $s$ , at least  $i$  elements of  $s$  remain uncovered. Thus, if the greedy algorithm were to choose  $s$  in that iteration, it would pay a cost per element of at most  $w_s/i$ . Thus, in this iteration, the greedy algorithm pays at most  $w_s/i$  per element covered. Thus, it charges element  $x_i$  at most  $w_s/i$  to be covered. Summing over  $i$ , the total amount charged to elements in  $s$  is at most  $w_s H_k$ . Summing over  $s \in C^*$  and noting that every element is in some set in  $C^*$ , the total amount charged to elements overall is at most  $\sum_{s \in C^*} w_s H_k = H_k \text{OPT}$ .  $\square$

The theorem was shown first for the unweighted case (each  $w_s = 1$ ) by Johnson [6], Lovász [9], and Stein [14], then extended to the weighted case by Chvátal [2].

Since then a few refinements and improvements have been shown, including the following:

**Theorem 2.** *Let  $\mathcal{S}$  be a set system over a universe with  $n$  elements and weights  $w_s \leq 1$ . The total weight of the cover  $C$  returned by the greedy algorithm is at most  $[1 + \ln(n/\text{OPT})]\text{OPT} + 1$  (compare to [13]).*

*Proof.* Assume without loss of generality that the algorithm covers the elements in order  $x_n, x_{n-1}, \dots, x_1$ . At the start of the iteration in which the algorithm covers  $x_i$ , there are at least  $i$  elements left to cover, and all of them could be covered using multiple sets of total cost  $\text{OPT}$ . Thus, there is some set that covers not-yet-covered elements at a cost of at most  $\text{OPT}/i$  per element.

Recall the charging scheme from the previous proof. By the preceding observation, element  $x_i$  is charged at most  $\text{OPT}/i$ . Thus, the total charge to elements  $x_n, \dots, x_i$  is at most  $(H_n - H_{i-1})\text{OPT}$ . Using the assumption that each  $w_s \leq 1$ , the charge to each of the remaining elements is at most 1 per element. Thus, the total charge to all elements is at most  $i - 1 + (H_n - H_{i-1})\text{OPT}$ . Taking  $i = 1 + \lceil \text{OPT} \rceil$ , the total charge is at most  $\lceil \text{OPT} \rceil + (H_n - H_{\lceil \text{OPT} \rceil})\text{OPT} \leq 1 + \text{OPT}(1 + \ln(n/\text{OPT}))$ .  $\square$

Each of the above proofs implicitly constructs a linear-programming primal-dual pair to show the approximation ratio. The same approximation ratios can be shown with respect to any fractional optimum (solution to the fractional set-cover linear program).

**Other results.** The greedy algorithm has been shown to have an approximation ratio of  $\ln n - \ln \ln n + O(1)$  [12]. For the special case of set systems whose duals have finite Vapnik-Chervonenkis (VC) dimension, other algorithms have substantially better approximation ratio [1]. Constant-factor approximation algorithms are known for geometric variants of the closely related  $k$ -median and facility location problems (see the K-median and Facility Location entry of this text).

The greedy algorithm generalizes naturally to many problems. For example, for minimizing a linear function subject to a submodular constraint (defined above), the natural extension of the greedy algorithm gives an  $H_k$ -approximate solution, where  $k = \max_{s \in \mathcal{S}} f(\{s\}) - f(\emptyset)$ , assuming  $f$  is integer-valued [10].

The set-cover problem generalizes to allow each element  $x$  to require an arbitrary number  $r_x$  of sets containing it to be in the cover. This generalization admits a polynomial-time  $O(\log n)$ -approximation algorithm [8].

The special case when each element belongs to at most  $r$  sets has a simple  $r$ -approximation algorithm [15, §15.2]. When the sets have uniform weights ( $w_s = 1$ ), the algorithm reduces to the following: select any maximal collection of elements, no two of which are contained in the same set; return all sets that contain a selected element.

The variant “Max k-coverage” asks for a set collection of total weight at most  $k$  covering as many of the elements as possible. This variant has a  $(1 - 1/e)$ -approximation algorithm [15, Problem 2.18] (see [7] for sets with non-uniform weights).

For a general discussion of greedy methods for approximate combinatorial optimization, see [5, Ch. 4].

Finally, under likely complexity-theoretic assumptions, the  $\ln n$  approximation ratio is essentially the best possible for any polynomial-time algorithm [3, 4].

### 3 APPLICATIONS

Set Cover and its generalizations and variants are fundamental problems with numerous applications. Examples include:

- selecting a small number of nodes in a network to store a file so that all nodes have a nearby copy,
- selecting a small number of sentences to be uttered to tune all features in a speech-recognition model [11],
- selecting a small number of telescope snapshots to be taken to capture light from all galaxies in the night sky,
- finding a short string having each string in a given set as a contiguous sub-string.

### 4 OPEN PROBLEMS [optional]

None to report.

### 5 EXPERIMENTAL RESULTS

None to report.

### 6 CROSS REFERENCES

#### EDITOR PLEASE FORMAT

Greedy Algorithms Entry 00423

Set Covering Entry 00275

K-median and Facility Location Entry 00479

## 7 RECOMMENDED READING

- [1] Hervé Brönnimann and Michael T. Goodrich, Almost optimal set covers in finite VC-dimension. *Discrete & Computational Geometry*, 14(4):463–479, 1995.
- [2] Vasek Chvátal. A greedy heuristic for the set-covering problem. *Mathematics of Operations Research*, 4(3):233–235, 1979.
- [3] Carsten Lund and Mihalis Yannakakis. On the hardness of approximating minimization problems. *Journal of the ACM*, 41(5):960–981, 1994.
- [4] Uriel Feige. A threshold of  $\ln n$  for approximating set cover. *Journal of the ACM*, 45(4):634–652, 1998.
- [5] Teo F. Gonzalez, *Handbook of Approximation Algorithms and Metaheuristics*. (Chapman & Hall/CRC Computer & Information Science Series), 2007.
- [6] David S. Johnson. Approximation algorithms for combinatorial problems. *Journal of Computer and System Sciences*, 9:256–278, 1974.
- [7] S. Khuller, A. Moss, and J. Naor. The budgeted maximum coverage problem. *Information Processing Letters*, 70(1):39–45, 1999.
- [8] Stavros G. Kolliopoulos and Neal E. Young. Tight approximation results for general covering integer programs. In *Proceedings of the forty-second annual IEEE Symposium on Foundations of Computer Science*, 522–528, 2001.
- [9] László Lovász. On the ratio of optimal integral and fractional covers. *Discrete Mathematics*, 13:383–390, 1975.
- [10] George L. Nemhauser and Laurence A. Wolsey. *Integer and Combinatorial Optimization*. John Wiley & Sons, New York, 1988.
- [11] J.P.H. van Santen and A.L. Buchsbaum. Methods for optimal text selection. In *Proceedings of the European Conference on Speech Communication and Technology (Rhodos, Greece)*, 2:553–556, 1997.
- [12] Petr Slavik. A tight analysis of the greedy algorithm for set cover. *Journal of Algorithms*, 25(2):237–254, 1997.
- [13] Aravind Srinivasan. Improved approximations of packing and covering problems. In *Proceedings of the twenty-seventh annual ACM Symposium on Theory of Computing*, 268–276, 1995.
- [14] S. K. Stein. Two combinatorial covering theorems. *J. Comb. Theory A*, 16:391–397, 1974.
- [15] V.V. Vazirani. *Approximation Algorithms*, Springer. 2001.