

More on "Set Cover Problem"

Problem Each set of weight (cost) 1.

Target: Find as few sets as possible.

Example

$$T_1 = \{1, 2, 3, 4, 5, 6\}$$

$$T_2 = \{5, 6, 8, 9\}$$

$$T_3 = \{1, 4, 7, 10\}$$

$$T_4 = \{2, 5, 7, 8, 11\}$$

$$T_5 = \{3, 6, 9, 12\}$$

$$T_6 = \{10, 11\}$$

$$T_1 \rightarrow T_4 \rightarrow T_5 \rightarrow T_3$$

$$C = \{T_1, T_3, T_4, T_5\}$$

(A solution but not an optimal one!)

Idea of greedy algorithm (Greedy-Set-Cover (X, F)).
mother set
collect

(*) At each stage, the greedy algorithm picks the set that covers the greatest number of elements not yet covered

1. $U \leftarrow X$ ($U = \{\text{uncovered elements}\}$)

2. $C \leftarrow \emptyset$ ($C = \{\text{covered elements}\}$)

3. While $U \neq \emptyset$

4. do select an $S \in F$ that maximizes $|S \cap U|$

5. $T_1 \leftarrow T_1 \cup S$

(2)

Theorem Let $d = \max\{|S| : S \in F \text{ and } F \text{ is a set cover of } S\}$.

Then Greedy-Set-Cover is a polynomial time $H(d)$ -approximation algorithm where $H(d) = \sum_{i=1}^d \frac{1}{i} = \ln d + O(1)$.

(Note: If $|C|$ and $|C^*|$ are the solution of Greedy-Set-Cover and optimal solution respectively, then $|C|/|C^*| \leq H(d)$.)

"Algorithm G"



Proof.

Assign a price "1" to each set $S \in F$ selected by the algorithm and distribute this price over the elements covered for the first time.

S_i : i^{th} subset selected

c_x : the price allocated to $x \in X$, that is covered for the first time at i^{th} iteration:

$$c_x = \frac{1}{|S_i \setminus (\bigcup_{j=1}^{i-1} S_j)|}$$

This implies that $|C| = \sum_{x \in X} c_x$. (C is a set cover of X)

Now, we use c_x where $x \in X$ to estimate $|C^*|$.

(*) The price assigned to the optimal cover is

$$\sum_{S \in C^*} \sum_{x \in S} c_x$$

(The sets in C^* may overlap.)

$$\Rightarrow \sum_{S \in C^*} \sum_{x \in S} c_x \geq \sum_{x \in X} c_x = |C|$$

Claim: $\sum_{x \in S} c_x \leq H(|S|)$.

Note that if this claim is true, then

$$|C| \leq \sum_{S \in C^*} H(|S|) \quad \left(\because |S| \leq \max\{|S| \mid S \in C^*\} \right)$$

$$\leq |C^*| \cdot H(d) \quad \text{where } d = \max\{|S| \mid S \in F\}$$

$$\Rightarrow \frac{|C|}{|C^*|} \leq H(d)$$

original collection

(At this stage, we don't know $\max\{|S| \mid S \in C^*\}$)

Let $u_i(S)$ be the number of ~~all~~ elements in S

remaining uncovered after S_1, S_2, \dots, S_i are selected to C .

$$u_i(S) = |C \setminus \bigcup_{j=1}^i S_j| \quad \dots \quad u_0(S) = |C|$$

(4)

Clearly $u_0(S) \geq u_1(S) \geq \dots \geq u_i(S)$.

Let k be the least index s.t. $u_k(S) = 0$.

Then $u_0(S) \geq u_1(S) \geq \dots \geq u_{k-1}(S) \geq u_k(S) = 0$

$$\Rightarrow \sum_{x \in S} c_x = \sum_{i=1}^k (u_{i-1}(S) - u_i(S)) \cdot \frac{1}{|S \setminus \bigcup_{j=1}^{i-1} S_j|}$$

Since $|S \setminus \bigcup_{j=1}^{i-1} S_j| \geq |S \setminus \bigcup_{j=1}^i S_j| = u_{i-1}(S)$,

$$\sum_{x \in S} c_x \leq \sum_{i=1}^k (u_{i-1}(S) - u_i(S)) \cdot \frac{1}{u_{i-1}(S)} \rightarrow u_{i-1}$$

$$= \sum_{i=1}^k (u_{i-1} - u_i) \frac{1}{u_{i-1}}$$

$$= \sum_{i=1}^k \left(1 - \frac{u_i}{u_{i-1}}\right)$$

$$\leq (H(u_0) - H(u_1)) + (H(u_1) - H(u_2)) + \dots + (H(u_{k-1}) - H(u_k))$$

$$= H(u_0) - H(0)$$

$$= H(|S|)$$



$$H(u_{i-1}) - H(u_i) = \frac{1}{u_{i-1}} + \frac{1}{u_{i-2}} + \dots + \frac{1}{u_i} \geq \frac{1}{u_{i-1}} \cdot (u_{i-1} - u_i)$$

Approximating Shortest Superstring via Set Cover

First, we give an example.

Let $S = \{ \text{CATGC} \textcircled{1}, \text{CTAAGT} \textcircled{2}, \text{GCTA} \textcircled{3}, \text{TTCA} \textcircled{4}, \text{ATGCATC} \textcircled{5} \}$.

$\mathcal{S} =$

$\left\{ \begin{array}{l} S_1 = \{ \text{CATGC} \textcircled{1} \} \\ S_2 = \{ \text{CTAAGT} \textcircled{2} \} \\ S_3 = \{ \text{GCTA} \textcircled{3} \} \\ S_4 = \{ \text{TTCA} \textcircled{4} \} \\ S_5 = \{ \text{ATGCATC} \textcircled{5} \} \end{array} \right.$	$S_6 = \{ \textcircled{1}, \textcircled{2} \}$ $S_7 = \{ \textcircled{1}, \textcircled{3} \}$ <i>bad!</i> $S_8 = \{ \textcircled{1}, \textcircled{4} \}$ <i>bad!</i> $S_9 = \{ \textcircled{1}, \textcircled{5} \}$ $S_{10} = \{ \textcircled{2}, \textcircled{3} \}$	$S_{11} = \{ \textcircled{2}, \textcircled{4} \}$ $S_{12} = \{ \textcircled{2}, \textcircled{5} \}$ $S_{13} = \{ \textcircled{3}, \textcircled{4} \}$ <i>bad!</i> $S_{14} = \{ \textcircled{3}, \textcircled{5} \}$ $S_{15} = \{ \textcircled{4}, \textcircled{5} \}$
--	---	---

Clearly, S can be covered by a sub-collection of \mathcal{S} .

But, we need a collection with minimum cost.

$c(S_1) = 5, c(S_2) = 6, c(S_3) = 4, c(S_4) = 4, c(S_5) = 7,$

$c(S_6) = 10 \quad \dots \quad c(S_{15}) = 10$

CATGC
 CTAAGT
 ↓
 CATGCTAAGT

TTCA
 ATGCATC
 ↓
 TTCATGCATC.

⑥

Definition:

Let $S = \{\vec{a}_1, \vec{a}_2, \dots, \vec{a}_n\}$. For strings \vec{a}_i and \vec{a}_j , if the last $k > 0$ symbols (characters) of \vec{a}_i are the same as the first k symbols of \vec{a}_j , let $\sigma_{i,j,k}$ denote the string obtained by overlapping these k symbols of \vec{a}_i and \vec{a}_j .

Let I be the set of $\sigma_{i,j,k}$'s for all valid choices of i, j, k , i.e., the set of all "good" superstrings of pairs of strings in S , i.e., $k > 0$.

(如果 $k=0$, 则不选!)

We use set $(\alpha_{i,j,k})$ to denote $\{\vec{a}_i, \vec{a}_j\}$.

Now, $F = \{\text{set}(\alpha) : \alpha \in SUI\}$.

(*) 如果能找到一个 set cover, 对应的 α 就可以把它们串起来成为 S 的一个 superstring (with \vec{a}_i 's as substrings).

⑥'

Algorithm G' , Greedy-Set-Cover with Cost (X, F)

1. $C \leftarrow \emptyset$

2. $U \leftarrow X$

3. While $U \neq \emptyset$ do

4. Find $S \in F \setminus C$ that minimizes $\alpha \stackrel{\text{def}}{=} \frac{\text{cost}(S)}{S \cap U}$.

5. for each $x \in S \cap U$ do

6. $\text{price}(x) \leftarrow \alpha$

7. $C \leftarrow C \cup \{S\}$

8. $U \leftarrow U \setminus S$

9. Return C

✓ Algorithm \tilde{S} (Superstring Algorithm)

1. Compute S.C. in Algorithm G'

2. Let $\{\text{set}(\sigma_1), \dots, \text{set}(\sigma_k)\}$ be the collection of sets returned by Algorithm G' .

3. return $\vec{\sigma} \stackrel{\text{def}}{=} \sigma_1 \sigma_2 \dots \sigma_k$.

②

Lemma Let OPT_{SC} denote the cost of an optimal solution to the S.C. instance in (X, F) , and OPT_{SS} denote the length of the shortest superstring of S . Then $OPT_{SC} \leq 2 \cdot OPT_{SS}$.

Proof. Let $set(a_1), \dots, set(a_x)$ be a solution S.C. (not optimal)

Let a_1, a_2, \dots, a_x be the corresponding strings. Since input string is covered by at most two a strings,

$$\sum_i |a_i| \leq 2 \cdot OPT_{SS}. \quad (?)$$

Theorem Greedy SCP $\leq 2 H_d \cdot OPT_{SSP}$