

# Lecture 4 (Part 2)

11

How good is a greedy algorithm?

(•) In general, it is OK to use, if no better algorithm exists.

We can solve the Set Cover Problem by using Linear Programming.

Definition (Linear Programming)

A "Linear Programming" is the problem of optimizing (minimizing or maximizing) a linear function subject to linear inequality constraints.

Example

Minimize

$$\sum_{j=1}^m c_j x_j$$

Subject to

$$\sum_{j=1}^n a_{ij} x_j \geq b_i, \quad i = 1, 2, \dots, m$$

$$x_j \geq 0$$

(Note) 由於答案不一定要 Integer, 有 polynomial time algorithm

"Integer programming is harder."

# Set Cover Problem

variable  $x_j$  for set  $S_j = \begin{cases} 1 & \text{if } S_j \text{ is selected;} \\ 0 & \text{if not.} \end{cases}$

上面例子

$S = \{1, 6\}$

$$\text{Min. } \sum_{j=1}^k c_j x_j = \text{OPT}$$

$x_j = \begin{cases} 1 & \text{if } S_j \text{ is selected} \\ 0 & \text{if not.} \end{cases}$

subject to

One of  $S_1$  and  $S_2$  must be chosen!

One of  $S_5, S_6$  and  $S_7$  must be chosen!

$$\left. \begin{aligned} x_1 + x_2 &\geq 1 \\ x_1 + x_3 + x_4 &\geq 1 \\ x_2 + x_3 + x_5 &\geq 1 \\ x_4 + x_6 + x_7 &\geq 1 \\ x_5 + x_6 + x_7 &\geq 1 \\ x_4 + x_5 + x_6 &\geq 1 \end{aligned} \right\} ?$$

Consider 1, 2, 3, 4, 5, 6  
Each element has to be "covered".

How to apply?

- 1 :  $S_1$  and  $S_2$
- 2 :  $S_1, S_3, S_4$
- 3 :  $S_2, S_3, S_5$
- 4 :  $S_4, S_6, S_7$
- 5 :  $S_5, S_6, S_7$
- 6 :  $S_4, S_5, S_6$

Here,  $c_j = |S_j|$ .

$x_i$ 's are in  $\{0, 1\}$ .

Review p. 8

- $S_1 = \{1, 2\}$
- $S_2 = \{1, 3\}$
- $S_3 = \{2, 3\}$
- $S_4 = \{2, 4, 6\}$
- $S_5 = \{3, 5, 6\}$
- $S_6 = \{4, 5, 6\}$
- $S_7 = \{4, 5\}$

More on "Set Cover Problem"

Problem Each set of weight (cost) 1.

Target: Find as few sets as possible.

Example cover  $[1, 12]$ .

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

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

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

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

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

$T_6 = \{10, 11\}$

先選一  $S_i$ ;  
再選尚未  
cover 的  
集合, 它增  
加最多的  
新元素!

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

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

(A solution but not  
an optimal one!)

Better

mother set collect

Idea of greedy algorithm (Greedy-Set-Cover  $(X, F)$ )

(\*) 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.  $U \leftarrow U \setminus S$ .
- Stop at  $U = \emptyset$

$|S \cap U|$

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

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  $C(T)$  and  $C(OPT)$  are the solution of Greedy-Set-Cover and optimal solution respectively, then  $\frac{C(T)}{C(OPT)} \leq H(d)$ .)

"Algorithm G"

Proof:

1 分给新加入的元素 (平均)

Assign a price "1" to each set  $S' \in \mathcal{F}$  selected by the algorithm and distribute this price over the elements covered for the first time. (evenly)

$S_i$ :  $i^{\text{th}}$  subset selected

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

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

$c_x = c_x^{(i)}$  for some  $i!$

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  $C$ .

$$u_i(S) = \left| S \setminus \bigcup_{j=1}^i S_j \right|, \quad u_0(S) = |S|.$$

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

$$\begin{aligned} \sum_{x \in S} c_x &\leq \sum_{i=1}^k (u_{i-1}(S) - u_i(S)) \cdot \frac{1}{u_{i-1}(S)} \xrightarrow{u_i} u_{i-1} \\ &= \sum_{i=1}^k (u_{i-1} - u_i) \frac{1}{u_{i-1}} \\ &= \sum_{i=1}^k \left(1 - \frac{u_i}{u_{i-1}}\right) \\ &\leq (H(u_0) - H(u_1)) + (H(u_1) - H(u_2)) + \dots + (H(u_{k-1}) - H(u_k)) \\ &= H(u_0) - H(0) \\ &= H(|S|)_0 \end{aligned}$$

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

Approximating Shortest Superstring via Set Cover

First, we give an example.

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

$S =$

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

Clearly,  $S$  can be covered by a sub-collection of  $S$ .

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

### 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  $k=0$ , 則不選!)

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

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

(\*) 如果能找到一個 set cover, 對應的  $\alpha$  就可以把它們串起來成為  $S$  的一個 superstring (with  $\vec{a}_i$ 's as substrings).



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

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

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