## LECTURE 21 (November 13th)

## More Approximation Algorithms

## Set Cover & Randomized Rounding

<u>Set Cover Problem</u> Let U be a universe of n elements

Let  $S_1, ..., S_m \subseteq U$  be a family of subsets of U with associated costs  $c_1, ..., c_m$ 

Goal: Pick a minimum cost set cover of U

L, collection of sets such that
there union equals U

Note This generalizes the vertex cover problem, since  $U = \{e_1, \dots e_m\}$  are the edges of the graph  $S_u = \{e_1 e_2 \text{ is incident on vertex } u\}$ 

Set Cover is also NP-hard, but we will see an approximation algorithm for a, with O (log n) approximation, using a LP relaxation.

# Integer Linear Programming Formulation

min 
$$\sum_{i=1}^{m} C_i \times_i$$

s.t.  $\sum_{i:u \in \mathcal{U}} x_i \ge 1$   $\forall u \in \mathcal{U}$  [every element is covered]

 $x_i \in \{0,1\}$   $i=1,...m$   $x_i = \begin{cases} 1 \text{ set } i \text{ is included} \end{cases}$ 

O set  $i$  is not included

#### LP relaxation\_

min 
$$\sum_{i=1}^{m} c_i x_i$$

s.t.  $\sum_{i=u \in u} x_i \ge 1$   $\forall u \in U$ 
 $0 \le x_i \le 1$   $i=1,...,m$ 

First, let us see what approximation we can obtain by using a deterministic rounding scheme analogous to vertex cover.

Let F be the maximum frequency of any element, i.e., maximum number of subsets any element appears in.

First we solve the LP to obtain a fractional solution x\*. Note that the LP objective value

$$OPT^* = \sum_{i=1}^m c_i x_i^*$$

satisfies that OPT & OPT where OPT is the cost of optimal set cover.

Then, we round it as follows

$$x_i = \begin{cases} 1 & \text{if } x_i \leq 1/F \\ 0 & \text{o/w} \end{cases}$$

Then,  $x = (x_1, \dots, x_m)$  is an integral solution that gives a set cover.

Moreover, cost of this set cover is

$$\left[\begin{array}{ccc} OPT & \leq \end{array}\right] \sum_{i=1}^{m} c_{i} \times_{i} & \leq F \sum_{i=1}^{m} c_{i} \times_{i} & = F \cdot OPT^{4}$$

Thus, this set cover is an F-approximation.

For vertex cover, F = 2, so we obtained a 2-approximation but F can be m in general, in which the approximation is trivial as the same can be achieved by including all subsets Si,..., Sm in the cover,

To obtaîn a much better approximation, we will use a randomized alporithm for rounding.

## Randomized Rounding for Set Cover

Solve the LP relaxation and obtain a fractional solution x as before.

- For each i=1,...m, round  $x_i^*\longrightarrow 1$  with probability  $x_i^*$ Li.e. include set  $S_i$  with probability  $x_i^*$ 2
- Repeat 2 until all elements are covered. 3

The intuition behind this is that the higher the xi value in the LP solution the higher probability of picking this set.

The above algorithm is harder to analyze so we consider a small variant:

- I Solve the LP relaxation and obtain a fractional solution x as before.
- Repeat log n + 2 times:

  For each i=1,...m, round x; 1 with probability x; 2 (i.e. include set S; with probability x; 2)
- 3 If the final integral solution does not cover all elements or cost is more than (4log n+8) factor of the LP solution, repeat [2]

To analyze this algorithm, let's see the cost of a single rounding step in 2

Let 
$$Y_i = \begin{cases} 1 & \text{if } S_i \text{ is picked} \\ 0 & \text{of } w \end{cases} \Rightarrow \text{This is a random variable}$$

After step 2 finishes,  $y = (y_1, ..., y_m)$  be the integral solution

Then 
$$\mathbb{E}\left[\sum_{i=1}^{m}c_{i},y_{i}\right] = \sum_{i=1}^{m}c_{i}$$
  $\mathbb{E}\left[y_{i}\right] = \sum_{i=1}^{m}c_{i},x_{i}^{*} = OPT^{*}$ 

So, the expected cost of the solution is exactly the LP objective value OPT

Over all the log n +2 iterations, the expected cost  $\leq$  (log n+2). OPT

By Markov's inequality, with probability 44, the cost of the final integral solution is < (410p n +8). OPT

What is the probability that this integral solution is not a set cover?

Consider any fixed element of the universe, say u,

IP [ u is not covered in any execution of the rounding step ]

= 
$$T$$
  $P[Si \text{ is not picked}]$ 

$$i: u \in Si$$

$$= T (1-x_i^*) \leq T e^{-x_i^*} = e^{-\sum_{i: u \in S_i} x_i^*} \leq \frac{1}{e}$$

$$i: u \in Si$$

If u is not covered in any of the log n+2 steps  $\int = \left(\frac{1}{e}\right)^{\log n+2} = \frac{1}{4n}$ 

By union bound

$$\mathbb{P}\left[\exists u \text{ that is not covered in any of the}\right] \leq n \cdot \frac{1}{4n} = \frac{1}{4}$$
 $\log n + 2$  steps

Thus, 
$$P \left[ \begin{array}{ccc} Final & integral & soln & after & step & 2 \\ is & a & cover & with & cast & \leq (4log n + 8) \cdot OPT \end{array} \right]$$

$$\leq \frac{1}{4} + \frac{1}{4} = \frac{1}{2}$$

Thus, in expectation, step 2 needs to be repeated 2 times and in the end we find a set cover whose cost is

Thus, we obtain a O((og n) approximation.

Note: It is NP-hard to obtain a better approximation of set cover.

## Hardness of Approximation

Unfortunately not all problems can be approximated beyond certain thresholds in poly-time. How do we prove that such problems are hard because these are not decision problems.

The basic idea is similar: reduce to a problem that is known to be NP-hard but one needs to take into account the approximation factors to convert it to a decision problem.

Let's see some examples.

#### Hardness of Traveling Salesman Problem

#### Traveling Salesman Problem

Given a list on n cities with distances d (i,j), find the shortest tour that visits each city exactly once and returns to the initial city.

We will prove the following

Theorem

For any function f(n) that can be computed in polynomial time in n, there is no polynomial-time f(n) approximation algorithm for the TSP on general weighted graphs unless P=NP.

Proof Sketch

(approximating)
If there is an algorithm for TSP, one can solve the
Hamiltonian Cycle problem in poly-calls to the TSP algorithm.

This is a decision problem:

Given a graph G=(V,E), is there a Hamiltonian Cycle in graph or not.

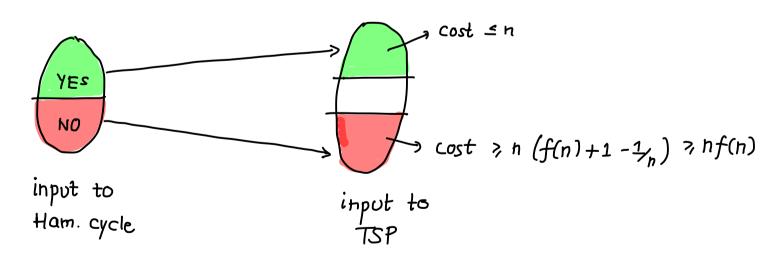
Since Hamiltonian Cycle is a known NP-hard problem, it follows that approximating TSP is NP-hard in general.

Reduction Given an instance  $G = (V_i E)$  for the Hamiltonian Cycle Problem we define a TSP instance as follows:

G' will be a complete graph on 
$$V \& d(i,j) = \begin{cases} 1 & \text{if } e \in E \\ n f(n) & \text{old} \end{cases}$$

(YES) If G had a Ham. cycle => G' has a tour with cost < n

(NO) If G didn't have a Ham. cycle  $\Rightarrow$  every tour in G' has cost  $\geqslant n f(n) + n - 1$ 



The main property of reductions that establish hardness is the gap between the two cases.

This proves that TSP is hard to approximate with any factor f(n).

How to deal with problems that are even hard to approximate, such as TSP?

Maybe our input have more structure that we are not using.

E.g. for TSP, our distances satisfy the triangle inequality in many cases of interest.

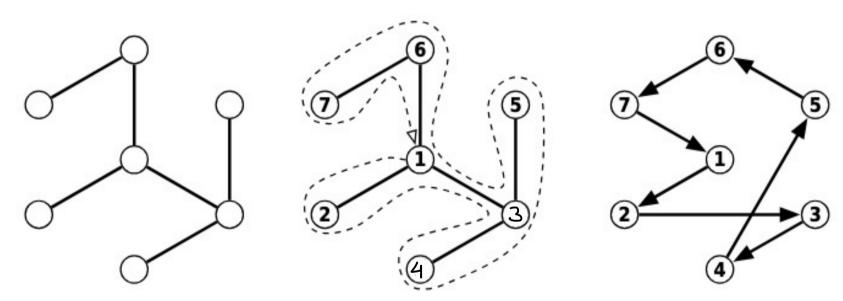
$$d(i,j) \leq d(i,k) + d(k,j)$$
  $\forall$  vertices  $i,j,k$ 

In this case, there is a simple 2-approximation algorithm for TSP.

In This is called the Metric TSP

#### Metric TSP algorithm

- 1 Compute a minimum spanning tree T of the weighted input graph
- 2 Perform a depth-first traversal of T numbering the vertices in this order
- B Return the tour obtained by visiting the vertices according to this numbering.



Theorem

This gives a 2-approximation to metric TSP.

This can be improved with new tools.

Proof First, consider the tour computed by walking the edges of MST in the order given by depth-first search. This is not a valid TSP tour since we will visit vertices more than once. But

cost of this "tour" = 2 · cost of MST, since each edge is traversed atmost twice

The final tour is obtained by removing duplicate vertices in the "tour" This does not increase the cost because of triangle inequality, going straight only costs less.

On the other hand, cost of MST < cost of optimal tour [Why?]

Thus, this gives us a 2-approximation