

In this lecture, a few ideas will come together: we will see how randomness, linear programming, and LP duality can be exploited to get good approximation algorithms for NP-complete problems. Specifically, we will see

- a LP-based approximation algorithm for WEIGHTED SET COVER.
- an approximation algorithm for WEIGHTED SET COVER that uses randomization to convert an LP to a good set cover
- an approximation algorithm for FACILITY LOCATION that simultaneously uses both a primal and the dual.

The algorithms in this lecture — and many more approximation algorithms — can be found in the book *Approximation Algorithms* by Vijay Vazirani.

1 Set Cover

Definition 1 WEIGHTED SET COVER: We are given a set X of n items, weights $w(x)$ for each $x \in X$, and m subsets $\mathcal{S} = \{S_1, S_2, \dots, S_m\}$ of these items. Goal: find the fewest number of these subsets needed to cover all the points. The decision problem also provides a number k and asks whether it is possible to cover all the points using k or fewer sets.

WEIGHTED SET COVER is NP-Complete. We can write it as an integer linear program. We have a binary variable x_S for each set:

$$\begin{aligned} & \text{minimize} && \sum_{S \in \mathcal{S}} w(S)x_S \\ & \text{subject to} && \sum_{S: e \in S} x_S \geq 1, \quad e \in X \\ & && x_S \in \{0, 1\}, \quad S \in \mathcal{S} \end{aligned} \tag{1}$$

The first set of constraints ensures that every item is in at least one chosen ($x_S = 1$) set. If we replace the (1) constraints with:

$$x_S \geq 0, \quad S \in \mathcal{S} \tag{2}$$

we obtain the *LP relaxation* of the integer program. The LP is easily solvable (using one of the algorithms we've already discussed). However, there is no guarantee that the relaxation will give an integer solution — instead it may return a *fractional* solution where some or all variables are real values rather than 0 or 1. If this happens, it's not clear immediately how we can extract a choice of sets from the LP solution, since we are not allowed to choose (say) 1/53rd of a set.

1.1 Deterministic Rounding of Weighted Set Cover (read on your own)

One way to deal with this is to try to *round* the fractional LP solution into an integer one. Suppose each item in X is in at most f sets. Then we can try the following algorithm:

1. Let x_S be an optimal solution to the LP relaxation.
2. Choose all sets S where $x_S \geq 1/f$

Theorem 2 *This gives a f -approximation algorithm for WEIGHTED SET COVER.*

Proof: Let $a \in X$ be any item. Since a is in $\leq f$ sets, at least one set S containing a must have $x_S \geq 1/f$. This means that every item must be covered by our choice of sets.

Our objective value increases by a factor $\leq f$ since at worst we change $x_S = 1/f \rightarrow x_S = 1$. So:

$$\text{rounded choice of sets} \leq f \times (\text{LP opt}) \leq f \times (\text{set cover opt})$$

■

This is a good approximation factor when f is small — that is, no item is in too many sets.

1.2 Randomized Rounding of Weighted Set Cover

We can get a good approximation algorithm by randomizing our procedure for converting the fractional solution to a choice of sets. In the LP relaxation solution, each x_S is in $[0, 1]$. That means we can think of x_S as a probability that x_S should be chosen. Consider the following *stage*:

- **Stage Algorithm:** For each $S \in \mathcal{S}$:
 - Add S to our cover C with probability x_S .

There are two things we need to ensure: (1) that there is a reasonably high probability that we get a feasible cover that covers all the elements, and (2) that the expected cost of the cover we get isn't too high.

For (1): Consider any element $a \in X$ that occurs in some number k sets. What is the chance that a is not covered after a single stage? The worst case is, for each S that contains a , the optimal LP solution gives $x_S = 1/k$. The probability that none of these is chosen:

$$\Pr[a \text{ is not covered}] \leq \left(1 - \frac{1}{k}\right)^k \leq \frac{1}{e}$$

To reduce this probability, we independently repeat multiple stages. In fact, we repeat the **Stage Algorithm** $c \ln n$ times, where we choose c so that:

$$\left(\frac{1}{e}\right)^{c \ln n} \leq \frac{1}{4n}$$

Using the union bound over all the items in X , we have:

$$\Pr[\text{some element is not covered}] \leq \frac{n}{4n} = 1/4.$$

For (2): The expected cost for one stage is $\sum_{S \in \mathcal{S}} x_S w(S)$, which is equal to the optimal value of the linear program OPT_{LP} . After $c \ln n$ stages, the expected cost is $OPT_{LP} \times c \ln n$. Recall Markov's inequality:

$$\Pr[X \geq b] \leq \frac{\mathbb{E}X}{b}$$

Applying this, we have

$$\Pr[\text{cost of our choice of sets} \geq 4 \times OPT_{LP} \times c \ln n] \leq \frac{1}{4}$$

Finally, we need both (1) and (2) to be true at the same time. This happens with at least probability $\geq 1/2$.

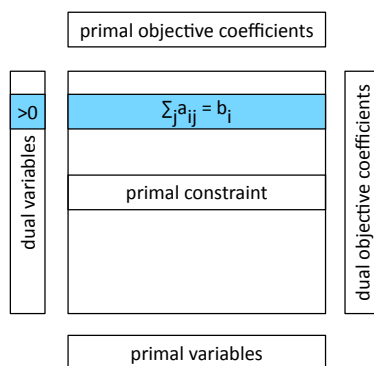
2 Complementary Slackness

Weak duality implies an interesting relationship between the non-zero variables in the primal and the constraints in the dual (and the non-zero dual variables and the primal constraints).

Theorem 3 (Complementary slackness conditions) *Let x and y be primal and dual feasible solutions for LPs written in standard form (using only \leq and \geq constraints). Let $\text{Constraint}(x_i)$ be the dual constraint corresponding to primal variable x_i and let $\text{Constraint}(y_i)$ be the primal constraint corresponding to dual variable y_i . Then x and y are optimal iff all the following conditions hold:*

- for all i , either $x_i = 0$ or $\text{Constraint}(x_i)$ is tight in the dual (holds with equality).
- for all j , either $y_j = 0$ or $\text{Constraint}(y_j)$ is tight in the primal.

Schematically:



Exercise! Prove the complementary slackness conditions from the fact that the when the primal and dual have finite optimal solutions, their values are equal. (Hint: see what weak duality implies in this case.)

3 Metric Facility Location

3.1 The problem

Here's the set up: Suppose you have a set of possible locations F at which to build fire stations. You also have locations C of buildings that need to be protected by fire stations. Buildings will be protected by the closest open fire station. Let $\{d_{ij}\}$ be the distances between these locations (both stations and buildings), and assume that d_{ij} form a metric, obeying the triangle inequality. It costs f_s to open fire station s .

Our goal: choose a subset I of fire stations F to minimize the following function:

$$\text{cost}(I) = \sum_{s \in I} f_s + \sum_{b \in B} \min_{s \in I} \{d_{bs}\}$$

In other words, we pay f_s for every station s we open, and we pay the distance between building b and it's closest open station. The first summation is called the *opening costs* and the second summation is called the *connection costs*.

3.2 The primal LP relaxation

We can encode this as an integer program, with two kinds of binary variables

- y_s is 1 iff station s is open (0 otherwise)
- x_{bs} is 1 iff building b is protected by station s (0 otherwise)

We can then encode the problem as an integer linear program (ILP):

$$\begin{array}{ll}
 \text{minimize} & \sum_{s \in F} f_s y_s + \sum_{s \in F, b \in C} d_{bs} x_{bs} \\
 \text{subject to} & \sum_{s \in F} x_{bs} \geq 1 & b \in C & \text{(covered)} \\
 & y_s - x_{bs} \geq 0 & b \in C, s \in F & \text{(open-iff-used)} \\
 & x_{bs} \in \{0, 1\} & & \text{(integral)} \\
 & y_s \in \{0, 1\} & & \text{(integral)}
 \end{array}$$

The constraints (covered) says that every building must be covered since it should be assigned to at least one station by setting x_{bs} to 1. The constraints (open-iff-used) require that x_{bs} can be 1 only when y_s is 1 — that is we can assign b to station s only when station s is open. Conversely, they require that if we assign b to station s (and $x_{bs} = 1$) then we must have opened station s .

We can construct the LP relaxation of this ILP by replacing the (integral) constraints by the requirement that the variables be ≥ 0 .

3.3 Dual of the LP relaxation

You can construct the dual entirely mechanically as we saw in a previous lecture using linear algebra. Here, to give a different perspective, we walk through the creation of the dual from the primal.

Recall for the dual we have variables for every constraint in the primal and constraints for every variable in the primal. Since we have two kinds of constraints ((covered) and (open-iff-used)) in the primal, we'll use 2 different symbols for the variables that they generate in the dual:

- α_b are variables corresponding to the (covered) constraints.
- β_{bs} are variables corresponding to the (open-iff-used) constraints.

Looking at the righthand side of the constraints in the primal gives us the coefficients of these variables in the dual objective:

$$\text{maximize} \quad \sum_{b \in C} \alpha_b$$

The coefficients of the β_{bs} constraints are 0, so they don't appear in the objective.

Each variable in the primal yields a constraint in the dual. The coefficient of y_s in the primal objective is f_s , so we have:

$$(\text{linear combination of dual variables corresponding to constraints containing } y_s) \leq f_s$$

The only constraints containing y_s are the “ β_{bs} ” constraints, and in each of those, y_s has coefficient 1. This gives us the constraint:

$$\sum_{b \in C} \beta_{bs} \leq f_s, \quad s \in F$$

in the dual.

The coefficient of the x_{bs} variable is d_{bs} in the primal objective. Further, x_{bs} appears in one (covered) constraint (the one for building b) and one (open-iff-used) constraints (the one for pair (b, s)). It has coefficient 1 in the first constraint and -1 in the second. This gives the following constraints in the dual:

$$\alpha_b - \beta_{bs} \leq d_{bs}, \quad s \in F, b \in C$$

Putting this together, we get the dual LP relaxation for metric facility location:

$$\text{maximize} \quad \sum_{b \in C} \alpha_b \tag{3}$$

$$\text{subject to} \quad \alpha_b - \beta_{bs} \leq d_{bs} \quad s \in F, b \in C \tag{4}$$

$$\sum_{b \in C} \beta_{bs} \leq f_s \quad s \in F \tag{5}$$

$$\alpha_b \geq 0 \tag{6}$$

$$\beta_{bs} \geq 0 \tag{7}$$

3.4 An interpretation of the dual

Suppose we think of the β_{bs} as payments from s toward opening s . The complimentary slackness condition on y_s says:

$$y_s > 0 \implies \sum_{b \in C} \beta_{bs} = f_s$$

In other words, if we open s , we have to fully pay for it in the dual. So our interpretation of β_{bs} makes some sense. Suppose we think of α_b as the total payment that building b makes. The complementary slackness condition for x_{bs} says:

$$x_{bs} > 0 \implies \alpha_b - \beta_{bs} = d_{bs}$$

In other words, if we connect b to s , our total payment $\alpha_b = d_{bs} + \beta_{bs}$. That is, b must pay for the road (d_{bs}) and its share of opening s .

3.5 A 4-approximation algorithm

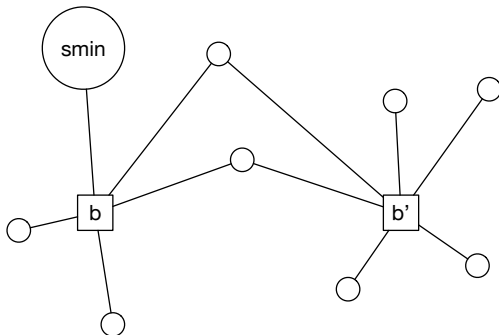
Consider the following algorithm:

1. Solve the primal and dual LP relaxations to obtain optimal $x_{bs}, y_s, \alpha_b, \beta_{bs}$.
2. Find the smallest α_b among the remaining buildings
 - (a) Define $N_b = \{s \mid x_{bs} > 0\}$ (the neighbor facilities of b)

- (b) Open the facility s_{\min} in N_b with the lowest $f_{s_{\min}}$ cost.
- (c) Assign (connect) every b' with $N_b \cap N_{b'} \neq \emptyset$ to s_{\min} .
- (d) Remove all assigned buildings and s_{\min}

3. Repeat step 2 until all facilities are assigned.

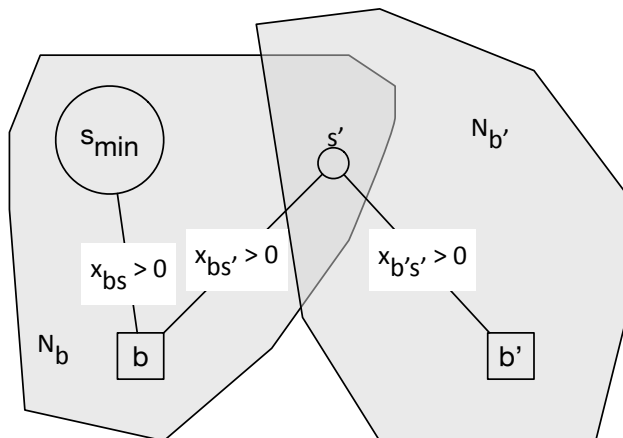
What's the intuition here? If the α_b is what b must pay, then starting with the cities that pay only a small amount makes sense. The strangest step is 2c. The idea here is that b' must be close to some neighbor of b , which is close to s_{\min} , so there is some hope that b' is close to s_{\min} because of the triangle inequality:



We now bound the connection costs and opening costs separately.

Theorem 4 *The total connection cost (sum of the chosen d_{bs}) of the above algorithm is $\leq 3 \sum_{b \in C} \alpha_b$.*

Proof: Consider when we looked at N_b and chose s_{\min} in step 2 of the algorithm. In step 2c, we connected b' to $s_{\min} = s$ too. So we had the following situation:



We know that each of the indicated x variables are non-zero since that is how the neighborhoods were defined. This means that the corresponding dual constraints are tight by complimentary slackness:

$$\begin{aligned} \alpha_b &= d_{bs} + \beta_{bs} \\ \alpha_b &= d_{bs'} + \beta_{bs'} \\ \alpha_{b'} &= d_{b's'} + \beta_{b's'} \end{aligned}$$

More importantly, this means:

$$\begin{aligned} d_{bs} &\leq \alpha_b \\ d_{bs'} &\leq \alpha_b \\ d_{b's'} &\leq \alpha_{b'} \end{aligned}$$

So:

$$\begin{aligned} d_{b's} &\leq d_{bs} + d_{bs'} + d_{b's'} && \text{(by triangle inequality)} \\ &\leq \alpha_b + \alpha_b + \alpha_{b'} && \text{(by above)} \\ &\leq 3\alpha_{b'} \end{aligned}$$

where the last inequality follows because we always chose the smallest α_b so $\alpha_b \leq \alpha_{b'}$.

So the total connection cost is at most $\sum_{b \in C} 3\alpha_b$, which is 3 times the optimal of the dual. ■

Theorem 5 *The total opening cost of the above algorithm is at most $OPT(P)$.*

Proof: Suppose b was chosen in step 2 of the algorithm and f_{\min} was the cost of the facility opened in N_b .

$$f_{\min} \leq f_{\min} \sum_{s \in N_b} x_{bs} \leq f_{\min} \sum_{s \in N_b} y_s \leq \sum_{s \in N_b} f_s y_s \quad (8)$$

The first inequality follows by the first set of primal constraints. The second inequality follows from the 2nd set of primal constraints. The final equality comes from the fact that we chose the cheapest facility.

All the considered N_b are disjoint, so the total opening cost is $\leq \sum_{s \in F} f_s y_s = OPT(P)$. ■

Putting it all together we have:

$$Our\ Algorithm \leq OPT(P) + 3OPT(D) = 4OPT(P) \leq 4OPT(int)$$

where $OPT(P)$, $OPT(D)$, and $OPT(int)$ are the optimal of the primal, dual, and integral solutions, respectively.