

LP-techniques for facility location and k -medians

Chs. 24, 25

Risto Hakala

`risto.m.hakala@tkk.fi`

April 10, 2008

- Facility location problem
 - A factor 3 approximation algorithm based on the primal-dual schema is presented.
- k -Median problem
 - A factor 6 approximation algorithm based on the previous algorithm is presented.

Facility location

Problem 24.1 (Metric uncapacitated facility location)

Let G be a complete bipartite graph with bipartition (F, C) , where F is a set of *facilities* and C is the set of *cities*. Let f_i be the cost of opening facility i , and c_{ij} be the cost of connecting city j to (opened) facility i . The connection costs satisfy the triangle inequality.

The problem is to find a subset $I \subseteq F$ of facilities that should be opened, and a function $\phi: C \rightarrow I$ assigning cities to open facilities in such a way that the cost of opening facilities and connecting cities is minimized.

- The problem is related to, e.g., locating proxy servers on the internet, and clustering.

Facility location problem as an integer program

- Let y_i and x_{ij} be indicator variables that denote whether facility i is open and whether city j is connected to facility i , respectively.
- We get the following IP:

$$\begin{array}{ll} \text{minimize} & \sum_{i \in F, j \in C} c_{ij} x_{ij} + \sum_{i \in F} f_i y_i \\ \text{subject to} & \sum_{i \in F} x_{ij} \geq 1, \quad j \in C \\ & y_i - x_{ij} \geq 0, \quad i \in F, j \in C \\ & x_{ij} \in \{0, 1\}, \quad i \in F, j \in C \\ & y_i \in \{0, 1\}, \quad i \in F \end{array}$$

LP-relaxation of the facility location problem

- As usual, the LP-relaxation is obtained by letting the domain of variables y_i and x_{ij} be $[0, \infty[$:

$$\text{minimize} \quad \sum_{i \in F, j \in C} c_{ij} x_{ij} + \sum_{i \in F} f_i y_i$$

$$\text{subject to} \quad \sum_{i \in F} x_{ij} \geq 1, \quad j \in C$$

$$y_i - x_{ij} \geq 0, \quad i \in F, j \in C$$

$$x_{ij} \geq 0, \quad i \in F, j \in C$$

$$y_i \geq 0, \quad i \in F$$

LP-relaxation of the facility location problem

- The dual program uses variables α_j and β_{ij} :

$$\begin{array}{ll} \text{maximize} & \sum_{j \in C} \alpha_j \\ \text{subject to} & \alpha_j - \beta_{ij} \leq c_{ij}, \quad i \in F, j \in C \\ & \sum_{j \in C} \beta_{ij} \leq f_i, \quad i \in F \\ & \alpha_j \geq 0, \quad j \in C \\ & \beta_{ij} \geq 0, \quad i \in F, j \in C \end{array}$$

- The variable β_{ij} can be viewed as the price paid by city j towards opening facility i .
- The variable α_j can be viewed as the total price paid by city j .

Primal-dual schema

- In the primal-dual schema, relaxed versions of complementary slackness conditions are used to guide the algorithm.
- The approximation factor is determined according to how much complementary slackness conditions have to be relaxed for them to be satisfied by the solution obtained from the algorithm.
- If a solution satisfies non-relaxed complementary slackness conditions, it is optimal.
- Hence, complementary slackness conditions define desirable properties for the algorithm.

Primal complementary slackness conditions

① $\forall i \in F, j \in C : x_{ij} > 0 \Rightarrow \alpha_j - \beta_{ij} = c_{ij}$

“The total price paid by the connected city goes towards making the connection and opening the facility.”

② $\forall i \in F : y_i > 0 \Rightarrow \sum_{j \in C} \beta_{ij} = f_i$

“Each open facility is fully paid for by the cities.”

Dual complementary slackness conditions

① $\forall j \in C : \alpha_j > 0 \Rightarrow \sum_{i \in F} x_{ij} = 1$

“All cities that pay anything must be connected to exactly one facility (with integral solutions).”

② $\forall i \in F, j \in C : \beta_{ij} > 0 \Rightarrow y_i = x_{ij}$

“A city does not contribute to opening any (open) facility besides the one that it is connected to.”

Primal-dual schema based algorithm

- The algorithm is divided into two parts: Phase 1 and Phase 2.
- Phase 1 finds a large dual feasible solution $(\vec{\alpha}, \vec{\beta})$ by changing only dual variables α_j and β_{ij} such that feasibility is maintained at all times.
- Phase 2 determines a primal (integral) feasible solution (\vec{x}, \vec{y}) based on the dual solution $(\vec{\alpha}, \vec{\beta})$.
- The approximation factor is determined by observing how much the complementary slackness conditions have to be relaxed in order for them to be satisfied.

Primal-dual schema based algorithm — Phase 1

Algorithm 24.2 — Phase 1

- Set $(\vec{\alpha}, \vec{\beta}) = (\vec{0}, \vec{0})$, time to 0, and define all cities to be *unconnected*.
- Do until all cities are *connected*:
 - Simultaneously raise α_j for each unconnected city j uniformly at unit rate, i.e., α_j grows 1 in unit time.
 - If $\alpha_j = c_{ij}$ for some edge (i, j) , declare this edge to be *tight* and start also raising β_{ij} uniformly at unit rate until j gets connected.
 - If $\sum_j \beta_{ij} = f_i$ for some facility i , declare this facility *temporarily open* and all unconnected cities having tight edges to i connected. Facility i is the *connecting witness* of cities that are connected to it.
 - If an unconnected city j gets a *tight* edge to a *temporarily open* facility, declare j connected.

Primal-dual schema based algorithm — Phase 1

After Phase 1,

- $\alpha_j - \beta_{ij} = c_{ij}$ for all tight edges (i, j) ,
- $\alpha_j < c_{ij}$ for all non-tight edges (i, j) ,
- $\sum_j \beta_{ij} = f_i$ for all temporarily open facilities i ,
- $\sum_j \beta_{ij} < f_i$ for all non-temporarily open facilities i .

Therefore, the fractional dual solution $(\vec{\alpha}, \vec{\beta})$ determined in Phase 1 is feasible.

Primal-dual schema based algorithm — Phase 2

- The set I of open facilities is picked from temporarily open facilities.
- Let
 - F_t denote the set of open facilities,
 - T denote the subgraph of G consisting of all “special” edges (i, j) such that $\beta_{ij} > 0$,
 - T^2 denote the graph that has edge (u, v) iff there is a path of length at most 2 between u and v in T , and
 - H denote the subgraph of T^2 induced on F_t .
- For city j , define $\mathcal{F}_j = \{i \in F_t \mid (i, j) \text{ is special}\}$.

Algorithm 24.2 — Phase 2

- Find any maximal independent set in H , say I .
- Iterate for all cities j :
 - If there is a facility $i \in \mathcal{F}_j$ that is opened ($i \in I$):
 - Set $\phi(j) = i$ and declare city j *directly connected*.
 - Else pick a tight edge (i', j) such that i' was the connecting witness for j .
 - If $i' \in I$, set $\phi(j) = i'$ and declare j *directly connected*.
 - If $i' \notin I$, pick a neighbor i of i' such that $i \in I$. Set $\phi(j) = i$ and declare j *indirectly connected*.
- Define a primal integral solution as follows:
 - Set $x_{ij} = 1$ iff $\phi(j) = i$.
 - Set $y_i = 1$ iff $i \in I$.

After Phase 2,

- there is a facility i such that $\phi(j) = i$ (i.e. $x_{ij} = 1$) for all cities j ,
- $\phi(j) = i$ (i.e. $x_{ij} = 1$) is set only whenever $i \in I$ (i.e. $y_i = 1$).

Therefore, the primal integral solution (\vec{x}, \vec{y}) determined in Phase 2 is feasible.

What about complementary slackness conditions?

Dual complementary slackness conditions

① $\forall j \in C : \alpha_j > 0 \Rightarrow \sum_{i \in F} x_{ij} = 1$

② $\forall i \in F, j \in C : \beta_{ij} > 0 \Rightarrow y_i = x_{ij}$

- Condition 1 is satisfied because $x_{ij} = 1$ is set for exactly one $i \in F$ for all $j \in C$.
- Condition 2 is satisfied because
 - $\phi(j) = i$ if $i \in \mathcal{F}_j$ is open, and
 - $\phi(j) \neq i$ if $i \in \mathcal{F}_j$ is not open.

What about complementary slackness conditions?

Primal complementary slackness conditions

- 1 $\forall i \in F, j \in C : x_{ij} > 0 \Rightarrow \alpha_j - \beta_{ij} = c_{ij}$
- 2 $\forall i \in F : y_i > 0 \Rightarrow \sum_{j \in C} \beta_{ij} = f_i$

- Condition 2 is satisfied because only temporarily opened facilities are opened fully.
- Condition 1 is satisfied for directly connected cities because a directly connected city j is connected to its facility i through a tight edge (i, j) .
- Condition 1 is not necessarily satisfied for indirectly connected cities since an indirectly connected city might not be connected to its facility through a tight edge.

What about complementary slackness conditions?

- In order to satisfy all conditions, the first primal complementary condition must be relaxed for indirectly connected cities j as follows:

$$(1/3)c_{\phi(j)j} \leq \alpha_j \leq c_{\phi(j)j}.$$

- This leads to an approximation algorithm that satisfies the inequality

$$\sum_{i \in F, j \in C} c_{ij}x_{ij} + 3 \sum_{i \in F} f_i y_i \leq 3 \sum_{j \in C} \alpha_j.$$

- Hence, the algorithm is a factor 3 approximation algorithm, but with a stronger inequality than typically.

Determination of the approximation factor

- Denote by α_j^f and α_j^e the contributions of city j to opening facilities and connection costs; $\alpha_j = \alpha_j^f + \alpha_j^e$.
- If j is indirectly connected, then $\alpha_j^f = 0$ and $\alpha_j^e = \alpha_j$.
- If j is directly connected, then $\alpha_j = c_{ij} + \beta_{ij}$, where $i = \phi(j)$.
- Let $\alpha_j^f = \beta_{ij}$ and $\alpha_j^e = c_{ij}$.

Determination of the approximation factor

Lemma 24.4

Let $i \in I$. Then,

$$\sum_{j:\phi(j)=i} \alpha_j^f = f_i.$$

Proof.

Since i is temporarily open at the end of Phase 1, it is completely paid for, i.e., $\sum_{j:\beta_{ij}>0} \beta_{ij} = f_i$. If city j has contributed to f_i , it must be directly connected to i . For each such city, $\alpha_j^f = \beta_{ij}$. Any other city j' that is connected to facility i must satisfy $\alpha_{j'}^f = 0$. The lemma follows. \square

Determination of the approximation factor

Corollary 24.5

$$\sum_{i \in I} f_i = \sum_{j \in C} \alpha_j^f.$$

Lemma 24.6

For an indirectly connected city j , $c_{ij} \leq 3\alpha_j^e$, where $i = \phi(j)$.

Determination of the approximation factor

Theorem 24.7

The primal and dual solutions constructed by the algorithm satisfy

$$\sum_{i \in F, j \in C} c_{ij} x_{ij} + 3 \sum_{i \in F} f_i y_i \leq 3 \sum_{j \in C} \alpha_j.$$

Proof.

For a directly connected city j , $c_{ij} = \alpha_j^e \leq 3\alpha_j^e$, where $\phi(j) = i$.
Combining with Lemma 24.6, we get

$$\sum_{i \in F, j \in C} c_{ij} x_{ij} \leq 3 \sum_{j \in C} \alpha_j^e.$$

Adding to this the equality stated in Corollary 24.5 multiplied by 3 gives the theorem. \square

Running time

- Denote $n_c = |C|$ and $n_f = |F|$.
- Sort all the edges by increasing cost — this gives the order and the times at which edges go tight.
- For each facility i , we maintain the number of cities that are currently contributing towards it, and the *anticipated time*, t_i , at which it would be completely paid for if no other event happens on the way.
- t_i 's are maintained in a binary heap so we can update each one and find the current minimum in $O(\log n_f)$ time.

- During the execution of the algorithm, t_i 's in the binary heap are updated whenever a facility is completely paid for or an edge goes tight.
- Each edge (i, j) will be considered at most twice: first, when it goes tight; second, when city j is declared connected.

Theorem 24.8

Algorithm 24.2 achieves an approximation factor of 3 for the facility location problem and has a running time of $O(m \log m)$, where $m = n_c \times n_f$ is the number of edges.

k -Median

Problem 24.1 (Metric k -Median)

Let G be a complete bipartite graph with bipartition (F, C) , where F is a set of *facilities* and C is the set of *cities*, and let k be a positive integer specifying the number of facilities that are allowed to be opened. Let c_{ij} be the cost of connecting city j to facility i . The connection costs satisfy the triangle inequality.

The problem is to find a subset $I \subseteq F, |I| \leq k$ of facilities that should be opened and a function $\phi: C \rightarrow I$ assigning cities to open facilities in such a way that the total connecting cost is minimized.

k -Median problem as an integer program

- Using indicator variables y_i and x_{ij} , we get the following IP:

$$\begin{array}{ll} \text{minimize} & \sum_{i \in F, j \in C} c_{ij} x_{ij} \\ \text{subject to} & \sum_{i \in F} x_{ij} \geq 1, \quad j \in C \\ & y_i - x_{ij} \geq 0, \quad i \in F, j \in C \\ & \sum_{i \in F} -y_i \geq -k \\ & x_{ij} \in \{0, 1\}, \quad i \in F, j \in C \\ & y_i \in \{0, 1\}, \quad i \in F \end{array}$$

LP-relaxation of the k -median problem

- The LP-relaxation is obtained by letting the domain of variables y_i and x_{ij} be $[0, \infty[$:

$$\text{minimize} \quad \sum_{i \in F, j \in C} c_{ij} x_{ij}$$

$$\text{subject to} \quad \sum_{i \in F} x_{ij} \geq 1, \quad j \in C$$

$$y_i - x_{ij} \geq 0, \quad i \in F, j \in C$$

$$\sum_{i \in F} -y_i \geq -k$$

$$x_{ij} \geq 0, \quad i \in F, j \in C$$

$$y_i \geq 0, \quad i \in F$$

LP-relaxation of the k -median problem

- Introducing the variables α_j and β_{ij} , we obtain the dual program:

$$\begin{array}{ll} \text{maximize} & \sum_{j \in C} \alpha_j - zk \\ \text{subject to} & \alpha_j - \beta_{ij} \leq c_{ij}, \quad i \in F, j \in C \\ & \sum_{j \in C} \beta_{ij} \leq f_i, \quad i \in F \\ & \alpha_j \geq 0, \quad j \in C \\ & \beta_{ij} \geq 0, \quad i \in F, j \in C \\ & z \geq 0 \end{array}$$

The high-level idea

- Consider a facility location problem, where the opening cost for each facility is $f_i = z$.
- By the strong duality theorem, the optimal fractional solutions (\vec{x}, \vec{y}) and $(\vec{\alpha}, \vec{\beta})$ satisfy

$$\sum_{i \in F, j \in C} c_{ij} x_{ij} + \sum_{i \in F} z y_i = \sum_{j \in C} \alpha_j.$$

- Suppose that the primal solution opens exactly k facilities, i.e., $\sum_i y_i = k$.

The high-level idea

- We obtain the equality

$$\sum_{i \in F, j \in C} c_{ij} x_{ij} = \sum_{j \in C} \alpha_j - zk.$$

- Hence, (\vec{x}, \vec{y}) and $(\vec{\alpha}, \vec{\beta}, z)$ are optimal fractional solutions to the k -median problem.
- Now, suppose we use Algorithm 24.2 to find primal integral and dual feasible solutions (\vec{x}, \vec{y}) and $(\vec{\alpha}, \vec{\beta})$ to the facility location problem such that exactly k facilities are opened.

The high-level idea

- By Theorem 24.7, the solutions satisfy

$$\sum_{i \in F, j \in C} c_{ij} x_{ij} + 3zk \leq 3 \sum_{j \in C} \alpha_j.$$

- Hence, (\vec{x}, \vec{y}) and $(\vec{\alpha}, \vec{\beta}, z)$ are primal integral and dual feasible solutions that satisfy

$$\sum_{i \in F, j \in C} c_{ij} x_{ij} \leq 3 \left(\sum_{j \in C} \alpha_j - zk \right).$$

- Algorithm 24.2 is a factor 3 approximation algorithm for the k -median problem *if* the value of z can be chosen such that exactly k facilities are opened.

The high-level idea

- It is not known how to choose z such that exactly k facilities are opened.
- To overcome this problem, the algorithm is used find solutions (\bar{x}^s, \bar{y}^s) and (\bar{x}^l, \bar{y}^l) to z_1 and z_2 , respectively, such that $k_1 < k$, $k_2 > k$, and $z_1 - z_2 \leq c_{\min}/(12n_c^2)$, where c_{\min} is the length of the shortest edge.
- The values of z_1 and z_2 are determined by conducting a binary search on the interval $[0, nc_{\max}]$, where n is the number of nodes and c_{\max} is the length of the longest edge.

The high-level idea

- The feasible (fractional) solution

$$(\vec{x}, \vec{y}) = a(\vec{x}^s, \vec{y}^s) + b(\vec{x}^l, \vec{y}^l), \quad ak_1 + bk_2 = k,$$

opens exactly k facilities. Here,

$$a = (k_2 - k)/(k_2 - k_1),$$

$$b = (k - k_1)/(k_2 - k_1).$$

Lemma 25.2

The cost of (\vec{x}, \vec{y}) is within a factor of $(3 + 1/n_c)$ of the cost of an optimal fractional solution to the k -median problem.

Randomized rounding

- An integral solution to the k -median problem is obtained from (\vec{x}, \vec{y}) using a randomized rounding procedure.
- Let A and B be the sets of opened facilities in solutions (\vec{x}^s, \vec{y}^s) and (\vec{x}^l, \vec{y}^l) , respectively.
- For each facility in A , find the closest facility in B , and form a set $B' \subset B$ using this facilities; if $|B'| < |A|$, arbitrarily include additional facilities from $B - B'$ into B' until $|B'| = |A| = k_1$.
- Open the facilities in A with probability $a = (k_2 - k)/(k_2 - k_1)$, and the facilities in B' with probability $b = (k - k_1)/(k_2 - k_1)$.
- Pick a set D of cardinality $k - k_1$ from $B - B'$, and open the facilities in it.

Randomized rounding

- The set of open facilities I is either $A \cup D$ or $B' \cup D$.
- Consider city j that is connected to facilities $i_1 \in A$ and $i_2 \in B$.
- If i_1 is open, set $\phi(j) = i_1$; if i_2 is open, set $\phi(j) = i_2$; otherwise, find the facility $i_3 \in B'$ that is closest to i_1 and set $\phi(j) = i_3$.
- Denote by $\text{cost}(j)$ the connection cost for city j in the fractional solution; $\text{cost}(j) = ac_{i_1j} + bc_{i_2j}$.

Lemma 25.3

The expected connection cost for city j in the integral solution, $\mathbf{E}[c_{\phi(j)j}]$, is $\leq (1 + \max(a, b))\text{cost}(j)$. Moreover, $\mathbf{E}[c_{\phi(j)j}]$ can be efficiently computed.

Lemma 25.4

Let (\vec{x}^k, \vec{y}^k) denote the integral solution obtained to the k -median problem by this randomized rounding procedure. Then,

$$\mathbf{E} \left[\sum_{i \in F, j \in C} c_{ij} x_{ij}^k \right] \leq (1 + \max(a, b)) \left(\sum_{i \in F, j \in C} c_{ij} x_{ij} \right),$$

and, moreover, the expected cost of the solution can be found efficiently.

Randomized rounding

- Derandomization is done by opening those sets which minimize the the previous expectation.
- The final approximation guarantee is $(1 + \max(a, b))(3 + 1/n_c) \leq (2 + 1/n_c)(3 + 1/n_c) < 6$.
- The binary search will make $O(L + \log n)$ probes, where $L = \log(c_{\max}/c_{\min})$.

Theorem 25.5

The algorithm achieves an approximation factor of 6 for the k -median problem, and has a running time of $O((m \log m)(L + \log n))$.

A Lagrangian relaxation technique for approximation algorithms

- A relaxation technique is a method in mathematical optimization for relaxing a strict requirement, e.g., by substituting it with another more easily handled requirement.
- Lagrangian relaxation technique consists of relaxing a (strict) constraint by moving it into the objective function, together with an associated Lagrangian multiplier λ .
- If the relaxed constrained is not satisfied, it induces a penalty on the objective function.

A Lagrangian relaxation technique for the k -median IP

- When applied to the k -median integer program, we obtain

$$\text{minimize} \quad \sum_{i \in F, j \in C} c_{ij} x_{ij} + \lambda \left(\sum_{i \in F} y_i - k \right)$$

$$\text{subject to} \quad \sum_{i \in F} x_{ij} \geq 1, \quad j \in C$$

$$y_i - x_{ij} \geq 0, \quad i \in F, j \in C$$

$$x_{ij} \in \{0, 1\}, \quad i \in F, j \in C$$

$$y_i \in \{0, 1\}, \quad i \in F$$

- This the facility location IP, where the cost of each facility has been set to λ , and an additional constant term $-\lambda k$ has been placed into the objective function.

Summary

- For the facility location problem, a factor 3 approximation algorithm based on the primal-dual schema was presented.
- This algorithm was used to construct a factor 6 approximation algorithm for the k -median problem.
- The primal-dual schema was used slightly differently here than in the previously presented algorithms.