

# An optimal bifactor approximation algorithm for the metric uncapacitated facility location problem

Jaroslav Byrka<sup>1</sup>

Centrum voor Wiskunde en Informatica,  
Kruislaan 413, NL-1098 SJ Amsterdam, Netherlands  
J.Byrka@cwi.nl

**Abstract.** We consider the metric uncapacitated facility location problem(UFL). In this paper we modify the  $(1+2/e)$ -approximation algorithm of Chudak and Shmoys to obtain a new  $(1.6774, 1.3738)$ -approximation algorithm for the UFL problem. Our linear programming rounding algorithm is the first one that touches the approximability limit curve  $(\gamma_f, 1 + 2e^{-\gamma_f})$  established by Jain et al. As a consequence, we obtain the first optimal approximation algorithm for instances dominated by connection costs.

Our new algorithm - when combined with a  $(1.11, 1.7764)$ -approximation algorithm proposed by Jain, Mahdian and Saberi, and later analyzed by Mahdian, Ye and Zhang - gives a 1.5-approximation algorithm for the metric UFL problem. This algorithm improves over the previously best known 1.52-approximation algorithm by Mahdian, Ye and Zhang, and it cuts the gap with the approximability lower bound by  $1/3$ .

Additionally, we show that a single randomized clustering procedure could be used instead of the greedy clustering used in the algorithms of Shmoys et al., Chudak et al., Sviridenko, and in the current paper.

---

<sup>1</sup> Supported by the EU Marie Curie Research Training Network ADONET, Contract No MRTN-CT-2003-504438

## 1 Introduction

The Uncapacitated Facility Location (UFL) problem is defined as follows. We are given a set  $\mathcal{F}$  of  $n_f$  facilities and a set  $\mathcal{C}$  of  $n_c$  clients. For every facility  $i \in \mathcal{F}$ , there is a nonnegative number  $f_i$  denoting the *opening cost* of the facility. Furthermore, for every client  $j \in \mathcal{C}$  and facility  $i \in \mathcal{F}$ , there is a *connection cost*  $c_{ij}$  between facility  $i$  and client  $j$ . The goal is to open a subset of the facilities  $\mathcal{F}' \subseteq \mathcal{F}$ , and connect each client to an open facility so that the total cost is minimized. The UFL problem is NP-complete, and max SNP-hard (see [4]). A UFL instance is *metric* if its *connection cost* function satisfies a kind of *triangle inequality*, namely if  $c_{ij} \leq c_{ij'} + c_{i'j} + c_{i'j'}$  for any  $i, i' \in \mathcal{C}$  and  $j, j' \in \mathcal{F}$ .

The UFL problem has a rich history starting in the 1960's. The first results on approximation algorithms are due to Cornuéjols, Fisher, and Nemhauser [1] who considered the problem with an objective function of maximizing the “profit” of connecting clients to facilities minus the cost of opening facilities. They showed that a greedy algorithm gives an approximation ratio of  $(1 - 1/e) = 0.632\dots$ , where  $e$  is the base of the natural logarithm. For the objective function of minimizing the sum of connection cost and opening cost, Hochbaum [2] presented a greedy algorithm with an  $O(\log n)$  approximation guarantee, where  $n$  is the number of clients. The first approximation algorithm with constant approximation ratio for the minimization problem where the connection costs satisfy the triangle inequality, was developed by Shmoys, Tardos, and Aardal [3]. Several approximation algorithms have been proposed for the metric UFL problem after that, see for instance [4–10]. Up to now, the best known approximation ratio was 1.52, obtained by Mahdian, Ye, and Zhang [10].

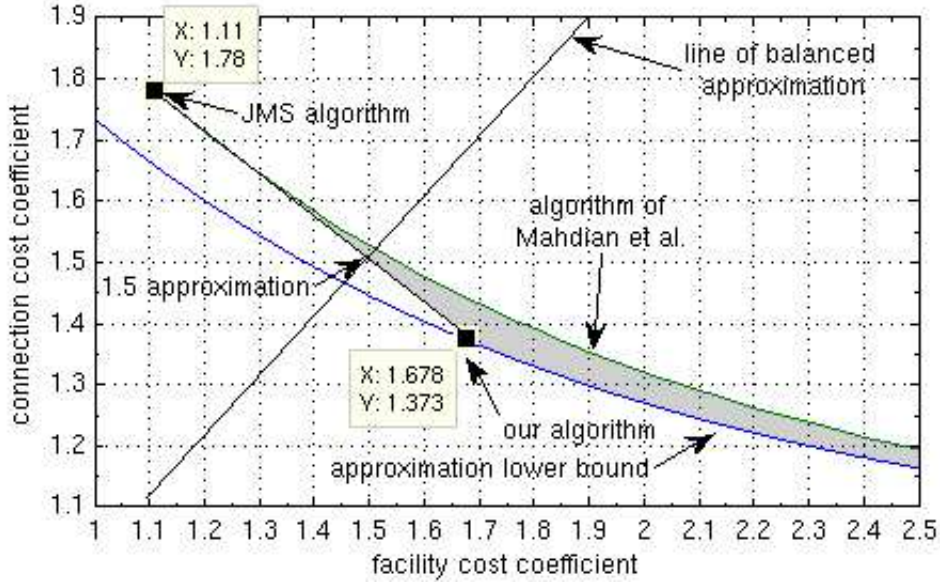
We will say that an algorithm is a  $\lambda$ -approximation algorithm for a minimization problem if it computes, in polynomial time, a solution that is at most  $\lambda$  times more expensive than the optimal solution. Specifically, for the UFL problem we define a notion of *bifactor approximation*. We say that an algorithm is a  $(\lambda_f, \lambda_c)$ -approximation algorithm if the solution it delivers has total cost at most  $\lambda_f \cdot F^* + \lambda_c \cdot C^*$ , where  $F^*$  and  $C^*$  denote, respectively, the facility and the connection cost of an optimal solution.

Guha and Khuller [4] proved by a reduction from Set Cover that there is no polynomial time  $\lambda$ -approximation algorithm for the metric UFL problem with  $\lambda < 1.463$ , unless  $NP \subseteq DTIME(n^{\log \log n})$ . Jain et al. [9] generalized this argument to show that the existence of a  $(\lambda_f, \lambda_c)$ -approximation algorithm with  $\lambda_c < 1 + 2e^{-\lambda_f}$  would imply  $NP \subseteq DTIME(n^{\log \log n})$ .

### 1.1 Our contribution

We modify the  $(1+2/e)$ -approximation algorithm of Chudak [6], see also Chudak and Shmoys [7], to obtain a new  $(1.6774, 1.3738)$ -approximation algorithm for the UFL problem. Our linear programming (LP) rounding algorithm is the first one that achieves an optimal bifactor approximation due to the matching lower bound of  $(\lambda_f, 1 + 2e^{-\lambda_f})$  established by Jain et al. In fact we obtain an algorithm for each point  $(\lambda_f, 1 + 2e^{-\lambda_f})$  such that  $\lambda_f \geq 1.6774$ , which means that we have an optimal approximation algorithm for instances dominated by connection cost (see Figure 1).

Our main technique is to sparsen the support graph corresponding to the LP solution before clustering. The motivation for this technique is the “irregularity” of instances that are potentially tight for the original algorithm of Chudak and Shmoys. We propose a way of measuring this irregularity and avoiding some of the too long links in the fractional solution to save on the connection cost of the final solution. In fact our clustering is the same as the one used by



**Fig. 1.** Bifactor approximation picture. The gray area corresponds to the improvement due to our algorithm.

Sviridenko in his 1.58-approximation algorithm [8], but we continue our algorithm in the spirit of Chudak and Shmoys' algorithm, which leads to a substantially easier analysis and an improved bifactor approximation guaranty.

Our new algorithm may be combined with the (1.11, 1.7764)-approximation algorithm of Jain et al. to obtain a 1.5-approximation algorithm for the UFL problem. This is an improvement over the previously best known 1.52-approximation algorithm of Mahdian et al., and it cuts of a 1/3 of the gap with the approximation lower bound by Guha and Khuler [4].

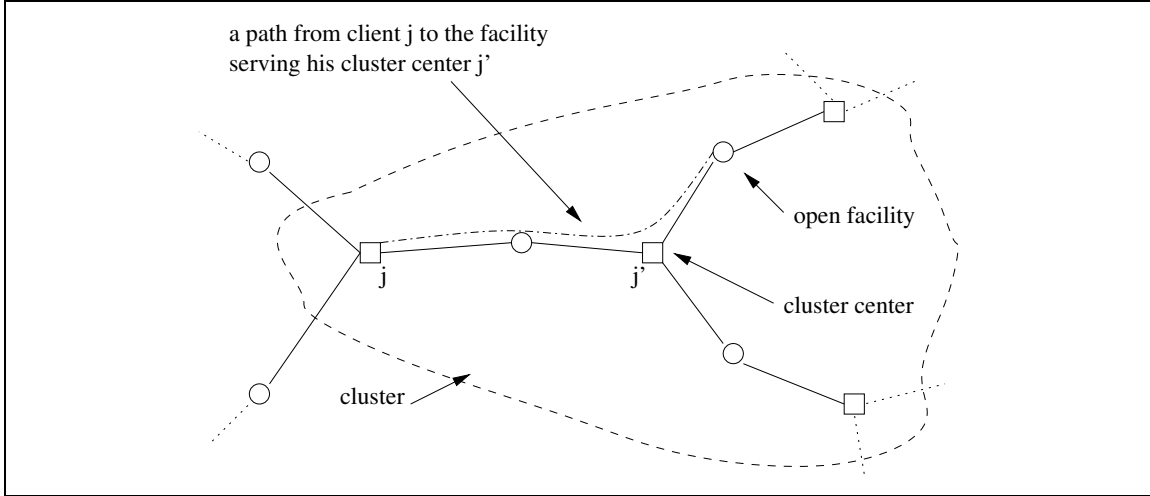
## 2 Preliminaries

We will review the concept of LP-rounding algorithms for the metric UFL problem. These are algorithms that first solve the linear relaxation of a given integer programming (IP) formulation of the problem, and then round the fractional solution to produce an integral solution with a value not too much higher than the starting fractional solution. Since the optimal fractional solution is at most as expensive as an optimal integral solution, we obtain an estimation of the approximation factor.

### 2.1 IP formulation and relaxation

The UFL problem has a natural formulation as the following integer programming problem.

$$\begin{aligned}
 & \text{minimize} && \sum_{i \in \mathcal{F}, j \in \mathcal{C}} c_{ij} x_{ij} + \sum_{i \in \mathcal{F}} f_i y_i \\
 & \text{subject to} && \sum_{i \in \mathcal{F}} x_{ij} = 1 && \text{for all } j \in \mathcal{C} && (1) \\
 & && x_{ij} - y_i \leq 0 && \text{for all } i \in \mathcal{F}, j \in \mathcal{C} && (2) \\
 & && x_{ij}, y_i \in \{0, 1\} && \text{for all } i \in \mathcal{F}, j \in \mathcal{C} && (3)
 \end{aligned}$$



**Fig. 2.** A cluster. If we make sure that at least one facility is open around a cluster center  $j'$ , then any other client  $j$  from the cluster may use this facility. Because the connection costs are assumed to be metric, the distance to this facility is at most the length of the shortest path from  $j$  to the open facility.

A linear relaxation of this IP formulation is obtained by replacing Condition (3) by the condition  $x_{ij} \geq 0$  for all  $i \in \mathcal{F}, j \in \mathcal{C}$ . The value of the solution to this LP relaxation will serve as a lower bound for the cost of the optimal solution. We will also make use of the following dual formulation of this LP.

$$\text{maximize } \sum_{j \in \mathcal{C}} v_j$$

$$\text{subject to } \sum_{j \in \mathcal{C}} w_{ij} \leq f_i \text{ for all } i \in \mathcal{F} \quad (4)$$

$$v_j - w_{ij} \leq c_{ij} \text{ for all } i \in \mathcal{F}, j \in \mathcal{C} \quad (5)$$

$$w_{ij} \geq 0 \text{ for all } i \in \mathcal{F}, j \in \mathcal{C} \quad (6)$$

## 2.2 Clustering

The first constant factor approximation algorithm for the metric UFL problem by Shmoys et al., but also the algorithms by Chudak and Shmoys, and by Sviridenko are based on the following clustering procedure. Suppose we are given an optimal solution to the LP relaxation of our problem. Consider the bipartite graph  $G$  with vertices being the facilities and the clients of the instance, and where there is an edge between a client  $j$  and a facility  $i$  if the corresponding variable  $x_{ij}$  in the optimal solution to the LP relaxation is positive. We call  $G$  a *support graph* of the LP solution. If two clients are both adjacent to the same facility in graph  $G$ , we will say that they are *neighbors* in  $G$ .

The clustering of this graph is a partitioning of clients into clusters together with a choice of a leading client for each of the clusters. This leading client is called a *cluster center*. Additionally we require that no two cluster centers are neighbors in the support graph. This property helps us to open one of the adjacent facilities for each cluster center. Formally we will say that a clustering is a function  $g : \mathcal{C} \rightarrow \mathcal{C}$  that assigns each client to the center of his cluster. For a picture of a cluster see Figure 2.

All the above mentioned algorithms use the following procedure to obtain the clustering. While not all the clients are clustered, choose greedily a new cluster center  $j$ , and build a cluster from  $j$  and all the neighbors of  $j$  that are not yet clustered. Obviously the outcome of this procedure is a proper clustering. Moreover, it has a desired property that clients are close to their cluster centers. Each of the mentioned LP-rounding algorithms uses a different greedy criterion for choosing new cluster centers. In our algorithm we will use the clustering with the greedy criterion of Sviridenko [8].

### 2.3 Scaling and greedy augmentation

The techniques described here are not directly used by our algorithm, but they help to explain why the algorithm of Chudak and Shmoys is close to optimal. We will discuss how scaling facility opening costs before running an algorithm, together with another technique called *greedy augmentation* may help to balance the analysis of an approximation algorithm for the UFL problem.

The greedy augmentation technique introduced by Guha and Khuller [4] (see also [5]) is the following. Consider an instance of the metric UFL problem and a feasible solution. For each facility  $i \in \mathcal{F}$  that is not opened in this solution, we may compute the impact of opening facility  $i$  on the total cost of the solution, also called the *gain* of opening  $i$ , denoted by  $g_i$ . The greedy augmentation procedure, while there is a facility  $i$  with positive gain  $g_i$ , opens a facility  $i_0$  that maximizes the ratio of saved cost to the facility opening cost  $\frac{g_i}{f_i}$ , and updates values of  $g_i$ . The procedure terminates when there is no facility whose opening would decrease the total cost.

Suppose we are given an approximation algorithm  $A$  for the metric UFL problem and a real number  $\delta \geq 1$ . Consider the following algorithm  $S_\delta(A)$ .

1. scale up all facility opening costs by a factor  $\delta$ ;
2. run algorithm  $A$  on the modified instance;
3. scale back the opening costs;
4. run the greedy augmentation procedure.

Following the analysis of Mahdian, Ye, and Zhang [10] one may prove the following lemma.

**Lemma 1.** *Suppose  $A$  is a  $(\lambda_f, \lambda_c)$ -approximation algorithm for the metric UFL problem, then  $S_\delta(A)$  is a  $(\lambda_f + \ln(\delta), 1 + \frac{\lambda_c - 1}{\delta})$ -approximation algorithm for this problem.*

This method may be applied to balance an  $(\lambda_f, \lambda_c)$ -approximation algorithm with  $\lambda_f \ll \lambda_c$ . However, our 1.5 approximation algorithm will be balanced differently. It will be a composition of two algorithms that have opposite imbalances.

### 3 Sparsening the graph of the fractional solution

Suppose that for a given UFL instance we have solved its LP relaxation, and that we have an optimal primal solution  $(x^*, y^*)$  and the corresponding optimal dual solution  $(v^*, w^*)$ . Such a fractional solution has facility cost  $F^* = \sum_{i \in \mathcal{F}} f_i y_i^*$  and connection cost  $C^* = \sum_{i \in \mathcal{F}, j \in \mathcal{C}} c_{ij} x_{ij}^*$ . Each client  $j$  has its share  $v_j$  of the total cost. This cost may again be divided into a client's fractional connection cost  $C_j^* = \sum_{i \in \mathcal{F}} c_{ij} x_{ij}^*$ , and his fractional facility cost  $F_j^* = v_j^* - C_j^*$ .

### 3.1 Motivation and intuition

The idea behind the sparsening technique is to make use of some irregularities of an instance if they occur. We call an instance *regular* if the facilities that fractionally serve a client  $j$  are all at the same distance from  $j$ . For such an instance the algorithm of Chudak and Shmoys produces a solution whose cost is bounded by  $F^* + (1 + \frac{2}{\epsilon})C^*$ , which also follows from our analysis in Section 4. It remains to use the technique described in section 2.3 to obtain an optimal 1.463...-approximation algorithm for such regular instances.

The instances that are not regular are called *irregular*. Difficult to understand are the irregular instances. In fractional solutions for these instances particular clients are fractionally served by facilities at different distances. Our approach is to divide facilities serving a client into two groups, namely *close* and *distant* facilities. We will remove links to distant facilities before the clustering step, so that if there are irregularities, distances to cluster centers should decrease.

We measure the local irregularity of an instance by comparing a fractional connection cost of a client to the average distance to his distant facilities. In the case of a regular instance, the sparsening technique gives the same results as technique described in section 2.3, but for irregular instances sparsening also takes some advantage of the irregularity.

### 3.2 Details

We will start by modifying the primal optimal fractional solution  $(x^*, y^*)$  by scaling the  $y$ -variables by a constant  $\gamma > 1$  to obtain a suboptimal fractional solution  $(x^*, \gamma \cdot y^*)$ . Now suppose that the  $y$ -variables are fixed, but that we now have a freedom to change the  $x$ -variables in order to minimize the total cost. For each client  $j$  we change the corresponding  $x$ -variables so that he uses his closest facilities in the following way. We choose an ordering of facilities with nondecreasing distances to client  $j$ . We connect client  $j$  to the first facilities in the ordering so that for any facilities  $i$  and  $i'$  such that  $i'$  is later in the ordering if  $x_{ij} < y_i$  than  $x_{i'j} = 0$ .

Without loss of generality, we may assume that this solution is complete (i.e. there are no  $i \in \mathcal{F}, j \in \mathcal{C}$  such that  $0 < x_{ij} < y_i$ ). Otherwise we may split facilities to obtain an equivalent instance with a complete solution - see [8][Lemma 1] for a more detailed argument.

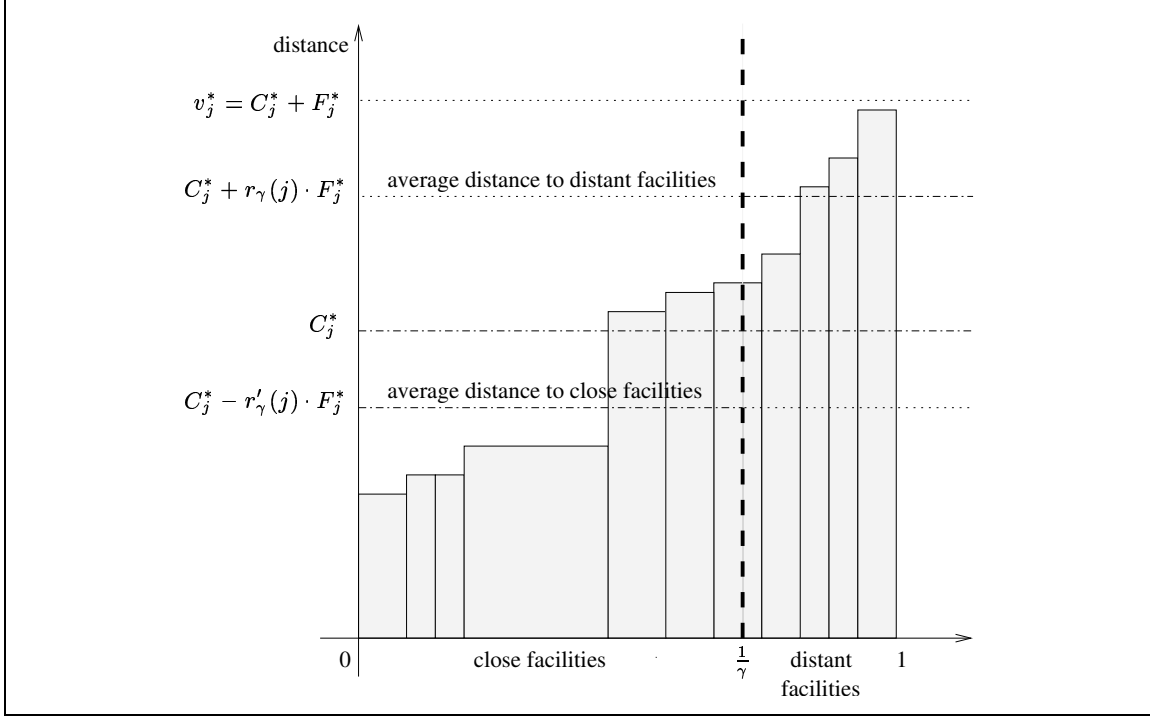
Let  $(\bar{x}, \bar{y})$  denote the obtained complete solution. For a client  $j$  we say that a facility  $i$  is one of *his close facilities* if it fractionally serves client  $j$  in  $(\bar{x}, \bar{y})$ . If  $\bar{x}_{ij} = 0$ , but facility  $i$  was serving client  $j$  in solution  $(x^*, y^*)$ , then we say, that  $i$  is a *distant* facility of client  $j$ .

**Definition 1.** Let  $r_\gamma(j) = \frac{\frac{\gamma}{\gamma-1} \sum_{i \in \{i \in \mathcal{F} | \bar{x}_{ij}=0\}} c_{ij} x_{ij}^* - C_j^*}{F_j^*}$  be the measure of a local irregularity of the instance. It is the average distance to a distant facility minus the fractional connection cost (which is the general average distance to both close and distant facilities), divided by the fractional facility cost of a client  $j$ .

Let  $r'_\gamma(j) = \frac{C_j^* - \sum_{i \in \mathcal{F}} c_{ij} \bar{x}_{ij}}{F_j^*} = r_\gamma(j) * (\gamma - 1)$  denote the fractional connection cost minus the average distance to a close facility, divided by the fractional facility cost of a client  $j$ .

Observe, that for every client  $j$  the following hold (see Figure 3):

- his average distance to a close facility equals  $D_{av}^C(j) = C_j^* - r'_\gamma(j) \cdot F_j^*$ ,
- his average distance to a distant facility equals  $D_{av}^D(j) = C_j^* + r_\gamma(j) \cdot F_j^*$ ,
- his maximal distance to a close facility is at most the average distance to a distant facility,  $D_{max}^C(j) \leq D_{av}^D(j) = C_j^* + r_\gamma(j) \cdot F_j^*$ .



**Fig. 3.** Distances to facilities serving client  $j$ ; the width of a rectangle corresponding to facility  $i$  is equal to  $x_{ij}^*$ . Figure explains the meaning of  $r_\gamma(j)$ .

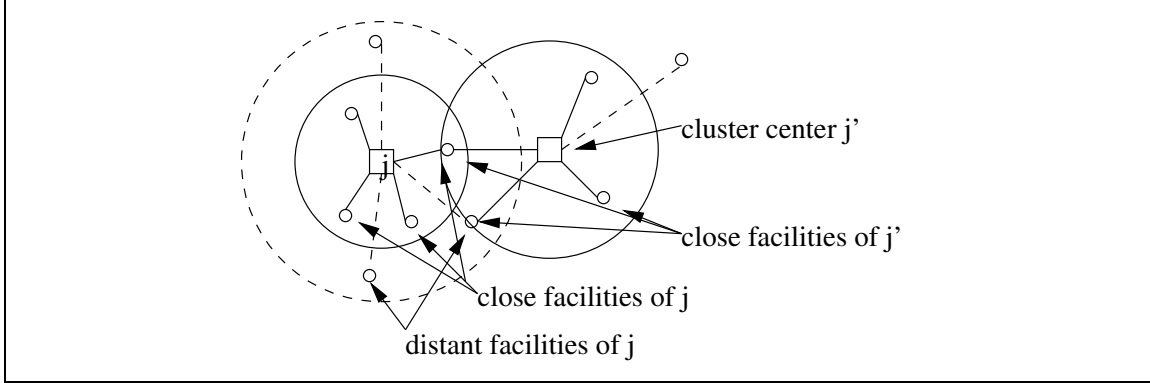
Consider a bipartite support graph  $G$  obtained from the solution  $(\bar{x}, \bar{y})$ , where each client is directly connected to his close facilities. We will greedily cluster this graph in each round choosing the cluster center to be an unclustered client  $j$  with the minimal value of  $D_{av}^C(j) + D_{max}^C(j)$ . With such a clustering, each cluster center has a minimal value of  $D_{av}^C(j) + D_{max}^C(j)$  among all clients in this cluster.

## 4 Our new algorithm

Consider the following algorithm  $A1(\gamma)$ :

1. Solve the LP relaxation of the problem to obtain a solution  $(x^*, y^*)$ .
2. Scale up the value of the facility opening variables  $y$  by a constant  $\gamma > 1$ , then change the value of the  $x$ -variables so as to use the closest possible fractionally open facilities (see Section 3.2).
3. If necessary, split facilities to obtain a complete solution  $(\bar{x}, \bar{y})$ .
4. Compute a greedy clustering for the solution  $(\bar{x}, \bar{y})$ , choosing as cluster centers unclustered clients minimizing  $D_{av}^C(j) + D_{max}^C(j)$ .
5. For every cluster center  $j$ , open one of his close facilities randomly with probabilities  $\bar{x}_{ij}$ .
6. For each facility  $i$  that is not a close facility of any cluster center, open it independently with probability  $\bar{y}_i$ .
7. Connect each client to an open facility that is closest to him.

In the analysis of this algorithm we will use the following result:



**Fig. 4.** Facilities that client  $j$  may consider: his close facilities, distant facilities, and close facilities of cluster center  $j'$ .

**Lemma 2.** Given  $n$  independent events  $e_1, e_2, \dots, e_n$  that occur with probabilities  $p_1, p_2, \dots, p_n$  respectively, the event  $e_1 \cup e_2 \cup \dots \cup e_n$  (i.e. at least one of  $e_i$ ) occurs with probability at least  $1 - \frac{1}{e^{\sum_{i=1}^n p_i}}$ , where  $e$  denotes the base of the natural logarithm.

**Theorem 1.** Algorithm A1( $\gamma = 1.67736$ ) produces a solution with expected cost  $E[\text{cost}(\text{SOL})] \leq 1.67736 \cdot F^* + 1.37374 \cdot C^*$

*Proof.* The expected facility opening cost of the solution is

$$E[F_{\text{SOL}}] = \sum_{i \in \mathcal{F}} f_i \bar{y}_i = \gamma \cdot \sum_{i \in \mathcal{F}} f_i y_i^* = \gamma \cdot F^*.$$

To bound the expected connection cost we show that for each client  $j$  there is an open facility within a certain distance with a certain probability. If  $j$  is a cluster center, one of his close facilities is open and the expected distance to this open facility is  $D_{av}^C(j) = C_j^* - r_\gamma(j) \cdot F_j^*$ .

If  $j$  is not a cluster center, he first considers his close facilities (see Figure 4). If any of them is open, the expected distance to the closest open facility is at most  $D_{av}^C(j)$ . From Lemma 2, with probability  $p_c \geq (1 - \frac{1}{e})$ , at least one close facility is open.

Suppose none of the close facilities of  $j$  is open, but at least one of his distant facilities is open. Let  $p_d$  denote the probability of this event. The expected distance to the closest facility is then at most  $D_{av}^D(j)$ .

If neither any close nor any distant facility of client  $j$  is open, then he connects himself to the facility serving his cluster center  $g(j) = j'$ . Again from Lemma 2, such an event happens with probability  $p_s \leq \frac{1}{e^\gamma}$ . In the following we will show that if  $\gamma < 2$  then the expected distance from  $j$  to the facility serving  $j'$  is at most  $D_{av}^D(j) + D_{max}^C(j') + D_{av}^C(j')$ . Let  $\mathcal{C}_j$  ( $\mathcal{D}_j$ ) be the set of close (distant) facilities of  $j$ . For any set of facilities  $X \subset \mathcal{F}$ , let  $d(j, X)$  denote the weighted average distance from  $j$  to  $i \in X$  (with values of opening variables  $y_i$  as weights).

If the distance between  $j$  and  $j'$  is at most  $D_{av}^D(j) + D_{av}^C(j')$ , then the remaining  $D_{max}^C(j')$  is enough for the distance from  $j'$  to any of his close facilities. Suppose now that the distance between  $j$  and  $j'$  is bigger than  $D_{av}^D(j) + D_{av}^C(j')$  (\*). We will bound  $d(j', \mathcal{C}_{j'} \setminus (\mathcal{C}_j \cup \mathcal{D}_j))$ , the average distance from cluster center  $j'$  to his close facilities that are neither close nor distant facilities of  $j$  (since the expected connection cost that we compute is on the condition that  $j$  was not served directly). The assumption(\*) implies that  $d(j', \mathcal{C}_j \cap \mathcal{C}_{j'}) > D_{av}^C(j')$ . Therefore, if  $d(j', \mathcal{D}_j \cap \mathcal{C}_{j'}) \geq D_{av}^C(j')$ , then  $d(j', \mathcal{D}_j \setminus (\mathcal{C}_j \cup \mathcal{D}_j)) \leq D_{av}^C(j')$  and the total distance from  $j$  is small enough.



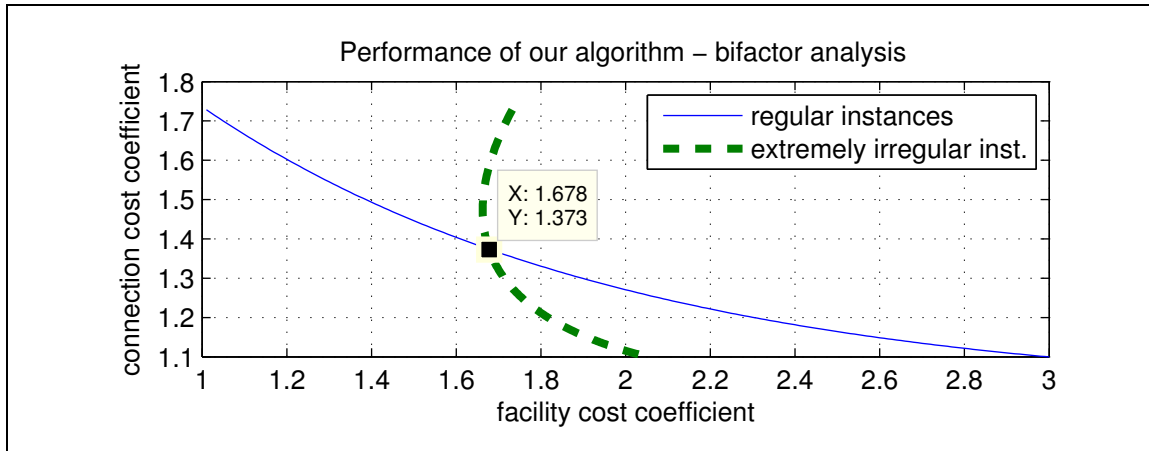
Suppose that  $d(j', \mathcal{D}_j \cap \mathcal{C}_{j'}) = D_{av}^C(j') - z$  for some positive  $z$  (\*\*). Let  $\hat{y} = \sum_{i \in (\mathcal{C}_{j'} \cup \mathcal{D}_j)} \bar{y}_i$  be the total fractional opening of facilities in  $\mathcal{C}_{j'} \cup \mathcal{D}_j$  in the modified fractional solution  $(\bar{x}, \bar{y})$ . From (\*) we conclude, that  $d(j, \mathcal{D}_j \cap \mathcal{C}_{j'}) \geq D_{av}^D(j) + z$ , which implies  $d(j, \mathcal{D}_j \setminus \mathcal{C}_{j'}) \leq D_{av}^D(j) - z \cdot \frac{\hat{y}}{\gamma - 1 - \hat{y}}$  (note that (\*\*) implies  $(\mathcal{D}_j \setminus \mathcal{C}_{j'}) \neq \emptyset$  and  $\gamma - 1 - \hat{y} > 0$ ), hence  $D_{max}^C(j) \leq D_{av}^D(j) - z \cdot \frac{\hat{y}}{\gamma - 1 - \hat{y}}$ . Combining this with assumption (\*) we conclude that the minimal distance from  $j'$  to a facility in  $\mathcal{C}_j \cap \mathcal{C}_{j'}$  is at least  $D_{av}^D(j) + D_{av}^C(j') - D_{max}^C(j) \geq D_{av}^C(j') + z \cdot \frac{\hat{y}}{\gamma - 1 - \hat{y}}$ . Assumption (\*\*) implies  $d(j', \mathcal{C}_{j'} \setminus \mathcal{D}_j) = D_{av}^C(j') + z \cdot \frac{\hat{y}}{1 - \hat{y}}$ . Concluding, if  $\gamma < 2$ , then  $d(j', \mathcal{C}_{j'} \setminus (\mathcal{D}_j \cup \mathcal{C}_j)) \leq D_{av}^C(j') + z \cdot \frac{\hat{y}}{\gamma - 1 - \hat{y}}$ . Therefore, the expected connection cost from  $j$  to a facility in  $\mathcal{C}_{j'} \setminus (\mathcal{D}_j \cup \mathcal{C}_j)$  is at most  $D_{max}^C(j) + D_{max}^C(j') + d(j', \mathcal{C}_{j'} \setminus (\mathcal{D}_j \cup \mathcal{C}_j)) \leq D_{av}^D(j) - z \cdot \frac{\hat{y}}{\gamma - 1 - \hat{y}} + D_{max}^C(j') + D_{av}^C(j') + z \cdot \frac{\hat{y}}{\gamma - 1 - \hat{y}} = D_{av}^D(j) + D_{max}^C(j') + D_{av}^D(j')$

Putting all the cases together, the expected total connection cost is

$$\begin{aligned}
E[C_{SOL}] &\leq \sum_{j \in \mathcal{C}} (p_c \cdot D_{av}^C(j) + p_d \cdot D_{av}^D(j) + p_s \cdot (D_{av}^D(j) + D_{max}^C(j') + D_{av}^C(j'))) \\
&\leq \sum_{j \in \mathcal{C}} ((p_c + p_s) \cdot D_{av}^C(j) + (p_d + 2p_s) \cdot D_{av}^D(j)) \\
&= \sum_{j \in \mathcal{C}} \left( (p_c + p_s) \cdot (C_j^* - r'_\gamma(j) \cdot F_j^*) + (p_d + 2p_s) \cdot (C_j^* + r_\gamma(j) \cdot F_j^*) \right) \\
&= ((p_c + p_d + p_s) + 2p_s) \cdot C^* \\
&\quad + \sum_{j \in \mathcal{C}} \left( (p_c + p_s) \cdot (-r_\gamma(j) \cdot (\gamma - 1) \cdot F_j^*) + (p_d + 2p_s) \cdot (r_\gamma(j) \cdot F_j^*) \right) \\
&= (1 + 2p_s) \cdot C^* + \sum_{j \in \mathcal{C}} \left( F_j^* \cdot r_\gamma(j) \cdot (p_d + 2p_s - (\gamma - 1) \cdot (p_c + p_s)) \right) \\
&\leq (1 + \frac{2}{e^\gamma}) \cdot C^* + \sum_{j \in \mathcal{C}} \left( F_j^* \cdot r_\gamma(j) \cdot (\frac{1}{e} + \frac{1}{e^\gamma} - (\gamma - 1) \cdot (1 - \frac{1}{e} + \frac{1}{e^\gamma})) \right)
\end{aligned}$$

Therefore, with  $\gamma = \gamma_0 \approx 1.67736$  such that  $\frac{1}{e} + \frac{1}{e^{\gamma_0}} - (\gamma_0 - 1) \cdot (1 - \frac{1}{e} + \frac{1}{e^{\gamma_0}}) = 0$ , we have  $E[C_{SOL}] \leq (1 + \frac{2}{e^{\gamma_0}}) \cdot C^* \leq 1.37374 \cdot C^*$ .  $\square$

The algorithm A1 with  $\gamma = 1 + \epsilon$  (for a sufficiently small positive  $\epsilon$ ) is essentially the algorithm of Chudak and Shmoys.



**Fig. 5.** Figure presents performance of our algorithm for different values of parameter  $\gamma$ . The solid line corresponds to regular instances with  $r_\gamma(j) = 0$  for all  $j$  and it coincides with the approximability lower bound curve. The dashed line corresponds to instances with  $r_\gamma(j) = 1$  for all  $j$ . For a particular choice of  $\gamma$  we get a horizontal segment connecting those two curves; for  $\gamma \approx 1.67736$  the segment becomes a single point. Observe that for instances dominated by connection cost only a regular instance may be tight for the analysis.

## 5 The 1.5-approximation algorithm

In this section we will combine our algorithm with an earlier algorithm of Jain et al. to obtain an 1.5-approximation algorithm for the metric UFL problem.

In 2002 Jain, Mahdian and Saberi [9] proposed a primal-dual approximation algorithm (the JMS algorithm). Using a dual fitting approach they have shown that it is a 1.61-approximation algorithm. In a later work of Mahdian, Ye and Zhang [10] the following was proven.

**Lemma 3** ([10]). *The cost of a solution produced by the JMS algorithm is at most  $1.11 \times F^* + 1.7764 \times C^*$ , where  $F^*$  and  $C^*$  are facility and connection costs in an optimal solution to the linear relaxation of the problem.*

**Theorem 2.** *Consider the solutions obtained with the A1 and JMS algorithms. The cheaper of them is expected to have a cost at most 1.5 times the cost of the optimal fractional solution.*

*Proof.* Consider the algorithm A2 that with probability  $p = 0.313$  runs the JMS algorithm and with probability  $1 - p$  runs the A1 algorithm. Suppose that you are given an instance, and  $F^*$  and  $C^*$  are facility and connection costs in an optimal solution to the linear relaxation of the problem for this instance. Consider the expected cost of the solution produced by algorithm A2 for this instance.  $E[\text{cost}] \leq p \cdot (1.11 \cdot F^* + 1.7764 \cdot C^*) + (1 - p) \cdot (1.67736 \cdot F^* + 1.37374 \cdot C^*) = 1.4998 \cdot F^* + 1.4998 \cdot C^* < 1.5 \cdot (F^* + C^*) \leq 1.5 \cdot \text{OPT}$ .  $\square$

Instead of the JMS algorithm we could take the algorithm of Machdian et al. [10] - the MYZ( $\delta$ ) algorithm that scales the facility costs by  $\delta$ , runs the JMS algorithms, scales back the facility costs and finally runs the greedy augmentation procedure. With a notation introduced in Section 2.3, the MYZ( $\delta$ ) algorithm is the  $S_\delta(\text{JMS})$  algorithm. The MYZ(1.504) algorithm was proven [10] to be a 1.52-approximation algorithm for the metric UFL problem. We may

change the value of  $\delta$  in the original analysis to observe that MYZ(1.1) is a (1.2053,1.7058)-approximation algorithm. This algorithm combined with our A1 (1.67736,1.37374)-approximation algorithm gives a 1.4991-approximation algorithm, which is even better than just using JMS and A1, but it gets more complicated and the additional improvement is tiny.

## 6 Universal randomized clustering procedure

In this section we discuss a different approach to clustering. We propose to modify the greedy clustering algorithm by choosing consecutive cluster centers randomly with uniform distribution. The output of such a process is obviously random, but we may still prove some statements about probabilities. The following lemma states that the obtained clustering is expected to be “fair”.

**Lemma 4.** *Given a graph  $G = (\mathcal{F} \cup \mathcal{C}, E)$  and assuming that a clustering  $g$  was obtained by the above described random process, for every two distinct clients  $j$  and  $j'$ , the probability that  $g(j) = j'$  is equal the probability that  $g(j') = j$ .*

*Proof.* Let  $C(G)$  denote the maximal (over the possible random choices of the algorithm) number of clusters that can be obtained from  $G$  with the random clustering procedure. The proof will be an induction on  $C(G)$ . Fix any  $j, j' \in \mathcal{C}$  such that  $j$  is a neighbor of  $j'$  in  $G$  (if they are not neighbors, neither  $g(j) = j'$  nor  $g(j') = j$  can occur). Suppose  $C(G) = 1$ , then  $Pr[g(j) = j'] = Pr[g(j') = j] = 1/|C|$ .

Let us now assume that  $C(G) > 1$ . There are two possibilities, either one of  $j, j'$  gets to the first cluster or they both avoid it. Consider the first case (the first chosen cluster center is either  $j$  or  $j'$  or one of their neighbors). If  $j$  ( $j'$ ) is chosen as a cluster center, then  $g(j') = j$  ( $g(j) = j'$ ). Since they are chosen with the same probability, the contribution of the first case to the probability of  $g(j') = j$  is equal to the contribution to the ppb. of  $g(j) = j'$ . If neither of them gets chosen as a cluster center but at least one gets into the new cluster, then neither  $g(j') = j$  nor  $g(j) = j'$  is possible.

Now consider the second case (neither of  $j$  and  $j'$  gets into the first cluster). Consider the graph  $G'$  obtained from  $G$  by removing the first cluster. The random clustering proceeds like it has just started with the graph  $G'$ , but the maximal number of possible clusters is smaller  $C(G') \leq C(G) - 1$ . Therefore, by the inductive hypothesis, in a random clustering of  $G'$  the ppb. that  $g(j') = j$  is equal the ppb. that  $g(j) = j'$ . Hence, the second case contribution to those probabilities for the clustering of the original graph  $G$  is also equal.  $\square$

If  $g(j) = j'$  in a clustering  $g$  of graph  $G$  we will say that client  $j'$  *offers a support* to client  $j$ . The main idea behind the clustering algorithms for the UFL problem is that we may afford to serve each cluster center directly (because they are never neighbors in  $G$ ) and all the other clients are offered a support from their cluster centers. A non-central client may either accept a support and connect himself via his cluster center (that is what all non-central clients do in the algorithm of Shmoys et al.), or he may try to get served locally, and if it fails, he will accept the support (this is the way the Chudak and Shmoys' algorithm works). In both those algorithms the probability that an offer of support is accepted is estimated to be constant. Therefore, we may modify those algorithms to use the random clustering procedure and do the following analysis.

For any two clients  $j$  and  $j'$ , the probability that  $j$  accepts a support of  $j'$  is equal to the probability that  $j'$  accepts the support of  $j$ . Let  $i$  be a facility on a shortest path from  $j$  to  $j'$ . When we compute the expected connection cost of a client  $j$ , we observe that with certain

probability  $p$  he accepts a support of  $j'$ . In such a case he must pay for the route via  $i$  and  $j'$  to the facility directly serving  $j'$ . In this situation we will say that  $j$  is paying only for the part until facility  $i$ , and the rest is paid by  $j'$ , but if  $j$  would be supporting  $j'$  he would have to pay a part of  $j'$ 's connection cost, which is the length of the path from  $i$  via  $j$  to the facility serving  $j$ . We may think of it as each client having a bank account, and when he accepts a support he makes a deposit, and when he offers a support and the support is accepted, then he withdraws money to pay a part of the connection cost of the supported client. From Lemma 4 we know that for a client  $j$  the probability that he will earn on  $j'$  is equal to the probability that he will lose on  $j'$ . Therefore, if the deposited amount is equal to the withdrawal, the expected net cash flow is zero.

The above analysis shows that randomizing the clustering phase of the above mentioned algorithms would not worsen their approximation ratios. Although it does not make much sense to use a randomized algorithm if it has no better performance guarantee, the random clustering has an advantage of allowing the analysis to be more local and uniform.

## 7 Concluding remarks

With the 1.52-approximation algorithm of Mahdian et al. it was not clear for the authors if a better analysis of the algorithm could close the gap with the approximation lower bound of 1.463 by Guha and Khuller. Byrka and Aardal [11] have recently given a negative answer to this question by constructing instances that are hard for the MYZ algorithm. Similarly, we now do not know if our new algorithm  $A1(\gamma)$  could be analyzed better to close the gap. Construction of hard instances for our algorithm remains an open problem.

The technique described in Section 2.3 enables to move the bifactor approximation guaranty of an algorithm along the approximability lower bound of Jain et al. (see Figure 1) towards higher facility opening costs. If we developed a technique to move the analysis in the opposite direction, together with our new algorithm, it would imply closing the approximability gap for the metric UFL problem. It seems that with such an approach we would have to face the difficulty of analyzing an algorithm that closes some of the previously opened facilities.

## References

1. G. Cornuéjols, M.L. Fisher, and G.L. Nemhauser. Location of bank accounts to optimize float: An analytic study of exact and approximate algorithms. *Management Science* **8**, pages 789–810, 1977.
2. D.S. Hochbaum. Heuristics for the fixed cost median problem. *Mathematical Programming* **22**, pages 148–162, 1982.
3. D. B. Shmoys, É. Tardos, and K. Aardal. Approximation algorithms for facility location problems (extended abstract). In *Proc. of the 29th ACM Symp. on Theory of Computing (STOC)*, pages 265–274, 1997.
4. S. Guha and S. Khuller. Greedy strikes back: Improved facility location algorithms. In *Proc. of the 9th ACM-SIAM Symp. on Discrete Algorithms (SODA)*, pages 228–248, 1998.
5. M. Charikar and S. Guha. Improved combinatorial algorithms for facility location and k-median problems. In *Proc. of the 40th IEEE Symp. on Foundations of Computer Science (FOCS)*, pages 378–388, 1999.
6. F. A. Chudak. Improved approximation algorithms for uncapacitated facility location. In *Proc. of the 6th Integer Programming and Combinatorial Optimization (IPCO)*, pages 180–194, 1998.
7. F. A. Chudak and D. B. Shmoys. Improved approximation algorithms for the uncapacitated facility location problem. *SIAM J. Comput.*, **33**(1), pages 1–25, 2003.
8. M. Sviridenko. An improved approximation algorithm for the metric uncapacitated facility location problem. In *Proc. of the 9th Integer Programming and Combinatorial Optimization (IPCO)*, pages 240–257, 2002.

9. K. Jain, M. Mahdian, and A. Saberi. A new greedy approach for facility location problems. In *Proc. of the 34th ACM Symp. on Theory of Computing (STOC)*, pages 731–740, 2002.
10. M. Mahdian, Y. Ye, and J. Zhang. Improved approximation algorithms for metric facility location problems. In *Proc. of the 5th International Workshop on Approximation Algorithms for Combinatorial Optimization Problems (APPROX)*, pages 229–242, 2002.
11. J. Byrka and K. Aardal. The approximation gap for the metric facility location problem is not yet closed. To appear in *Operations Research Letters (ORL)*