

# LP-rounding Algorithms for the Fault-Tolerant Facility Placement Problem

Li Yan and Marek Chrobak  
Department of Computer Science  
University of California at Riverside

## Abstract

The Fault-Tolerant Facility Placement problem (FTFP) is a generalization of the classic Uncapacitated Facility Location Problem (UFL). In FTFP we are given a set of facility sites and a set of clients. Opening a facility at site  $i$  costs  $f_i$  and connecting client  $j$  to a facility at site  $i$  costs  $d_{ij}$ . We assume that the connection costs (distances)  $d_{ij}$  satisfy the triangle inequality. Multiple facilities can be opened at any site. Each client  $j$  has a demand  $r_j$ , which means that it needs to be connected to  $r_j$  different facilities (some of which could be located on the same site). The goal is to minimize the sum of facility opening cost and connection cost.

The main result of this paper is a 1.575-approximation algorithm for FTFP, based on LP-rounding. The algorithm first reduces the demands to values polynomial in the number of sites. Then it uses a technique that we call adaptive partitioning, which partitions the instance by splitting clients into unit demands and creating a number of (not yet opened) facilities at each site. It also partitions the optimal fractional solution to produce a fractional solution for this new instance. The partitioned instance satisfies a number of properties that allow us to exploit existing LP-rounding methods for UFL to round our partitioned solution to an integral solution, preserving the approximation ratio. In particular, our 1.575-approximation algorithm is based on the ideas from the 1.575-approximation algorithm for UFL by Byrka *et al.*, with changes necessary to satisfy the fault-tolerance requirement.

# 1 Introduction

In the *Fault-Tolerant Facility Placement* problem (FTFP), we are given a set  $\mathbb{F}$  of *sites* at which facilities can be built, and a set  $\mathbb{C}$  of *clients* with some demands that need to be satisfied by different facilities. A client  $j \in \mathbb{C}$  has demand  $r_j$ . Building one facility at a site  $i \in \mathbb{F}$  incurs a cost  $f_i$ , and connecting one unit of demand from client  $j$  to a facility at site  $i$  costs  $d_{ij}$ . Throughout the paper we assume that the connection costs (distances)  $d_{ij}$  form a metric, that is, they are symmetric and satisfy the triangle inequality. In a feasible solution, some number of facilities, possibly zero, are opened at each site  $i$ , and demands from each client are connected to those open facilities, with the constraint that demands from the same client have to be connected to different facilities. Note that any two facilities at the same site are considered different.

It is easy to see that if all  $r_j = 1$  then FTFP reduces to the classic Uncapacitated Facility Location problem (UFL). If we add a constraint that each site can have at most one facility built on it, then the problem becomes equivalent to the Fault-Tolerant Facility Location problem (FTFL). One implication of the one-facility-per-site restriction in FTFL is that  $\max_{j \in \mathbb{C}} r_j \leq |\mathbb{F}|$ , while in FTFP the values of  $r_j$ 's can be much bigger than  $|\mathbb{F}|$ .

The UFL problem has a long history; in particular, great progress has been achieved in the past two decades in developing techniques for designing constant-ratio approximation algorithms for UFL. Shmoys, Tardos and Aardal [16] proposed an approach based on LP-rounding, that they used to achieve a ratio of 3.16. This was then improved by Chudak [5] to 1.736, and later by Sviridenko [17] to 1.582. The best known “pure” LP-rounding algorithm is due to Byrka *et al.* [3] with ratio 1.575. Byrka and Aardal [2] gave a hybrid algorithm that combines LP-rounding and dual-fitting (based on [10]), achieving a ratio of 1.5. Recently, Li [13] showed that, with a more refined analysis and randomizing the scaling parameter used in [2], the ratio can be improved to 1.488. This is the best known approximation result for UFL. Other techniques include the primal-dual algorithm with ratio 3 by Jain and Vazirani [11], the dual fitting method by Jain *et al.* [10] that gives ratio 1.61, and a local search heuristic by Arya *et al.* [1] with approximation ratio 3. On the hardness side, UFL is easily shown to be NP-hard, and it is known that it is not possible to approximate UFL in polynomial time with ratio less than 1.463, provided that  $\text{NP} \not\subseteq \text{DTIME}(n^{O(\log \log n)})$  [6]. An observation by Sviridenko strengthened the underlying assumption to  $\text{P} \neq \text{NP}$  (see [19]).

FTFL was first introduced by Jain and Vazirani [12] and they adapted their primal-dual algorithm for UFL to obtain a ratio of  $3 \ln(\max_{j \in \mathbb{C}} r_j)$ . All subsequently discovered constant-ratio approximation algorithms use variations of LP-rounding. The first such algorithm, by Guha *et al.* [7], adapted the approach for UFL from [16]. Swamy and Shmoys [18] improved the ratio to 2.076 using the idea of pipage rounding introduced in [17]. Most recently, Byrka *et al.* [4] improved the ratio to 1.7245 using dependent rounding and laminar clustering.

FTFP is a natural generalization of UFL. It was first studied by Xu and Shen [20], who extended the dual-fitting algorithm from [10] to give an approximation algorithm with a ratio claimed to be 1.861. However their algorithm runs in polynomial time only if  $\max_{j \in \mathbb{C}} r_j$  is polynomial in  $O(|\mathbb{F}| \cdot |\mathbb{C}|)$  and the analysis of the performance guarantee in [20] is flawed<sup>1</sup>. To date, the best approximation ratio for FTFP in the literature is 4, established by Yan and Chrobak [21], while the only known lower bound is the 1.463 lower bound for UFL from [6], as UFL is a special case of FTFP. If all demand values  $r_j$  are equal, the problem can be solved by simple scaling and applying LP-rounding

---

<sup>1</sup>Confirmed through private communication with the authors.

algorithms for UFL. This does not affect the approximation ratio, thus achieving ratio 1.575 for this special case (see also [14]).

The main result of this paper is an LP-rounding algorithm for FTFP with approximation ratio 1.575, matching the best ratio for UFL achieved via the LP-rounding method [3] and significantly improving our earlier bound in [21]. In Section 3 we prove that, for the purpose of LP-based approximations, the general FTFP problem can be reduced to the restricted version where all demand values are polynomial in the number of sites. This *demand reduction* trick itself gives us a ratio of 1.7245, since we can then treat an instance of FTFP as an instance of FTFL, by creating a sufficient (but polynomial) number of facilities at each site and using the algorithm from [4].

The reduction to polynomial demands suggests an approach where clients' demands are split into unit demands. These unit demands can be thought of as "unit-demand clients", and a natural approach would be to adapt LP-rounding methods from [8, 5, 3] to this new set of unit-demand clients. Roughly, these algorithms iteratively pick a client that minimizes a certain cost function (that varies for different algorithms) and open one facility in the neighborhood of this client. The remaining clients are then connected to these open facilities. In order for this to work, we also need to convert the optimal fractional solution  $(\mathbf{x}^*, \mathbf{y}^*)$  of the original instance into a solution  $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$  of the modified instance which then can be used in the LP-rounding process. This can be thought of as partitioning the fractional solution, as each connection value  $x_{ij}^*$  must be somehow divided between the  $r_j$  unit demands of client  $j$ . In Section 4 we formulate a set of properties required for this partitioning to work. For example, one property guarantees that we can connect demands to facilities so that two demands from the same client are connected to different facilities. Then we present our *adaptive partitioning* technique that computes a partitioning with all the desired properties. Using adaptive partitioning we were able to extend the algorithms for UFL from [8, 5, 3] to FTFP. We illustrate the fundamental ideas of our approach in Section 5, showing how they can be used to design an LP-rounding algorithm with ratio 3. In Section 6 we refine the algorithm to improve the approximation ratio to  $1 + 2/e \approx 1.736$ . Finally, in Section 7, we improve it even further to 1.575 – the main result of this paper.

Summarizing, our contributions are two-fold: One, we show that the existing LP-rounding algorithms for UFL can be extended to a much more general problem FTFP, retaining the approximation ratio. We believe that, should even better LP-rounding algorithms be developed for UFL in the future, using our demand reduction and adaptive partitioning methods, it should be possible to extend them to FTFP. In fact, some improvement of the ratio should be achieved by randomizing the scaling parameter  $\gamma$  used in our algorithm, as Li showed in [13] for UFL. (Since the ratio 1.488 for UFL in [13] uses also dual-fitting algorithms [15], we would not obtain the same ratio for FTFP yet using only LP-rounding.)

Two, our ratio of 1.575 is significantly better than the best currently known ratio of 1.7245 for the closely-related FTFL problem. This suggests that in the fault-tolerant scenario the capability of creating additional copies of facilities on the existing sites makes the problem easier from the point of view of approximation.

## 2 The LP Formulation

The FTFP problem has a natural Integer Programming (IP) formulation. Let  $y_i$  represent the number of facilities built at site  $i$  and let  $x_{ij}$  represent the number of connections from client  $j$  to facilities at site  $i$ . If we relax the integrality constraints, we obtain the following LP:

$$\begin{aligned}
& \text{minimize} && \text{cost}(\mathbf{x}, \mathbf{y}) = \sum_{i \in \mathbb{F}} f_i y_i + \sum_{i \in \mathbb{F}, j \in \mathbb{C}} d_{ij} x_{ij} \\
& \text{subject to} && y_i - x_{ij} \geq 0 \quad \forall i \in \mathbb{F}, j \in \mathbb{C} \\
& && \sum_{i \in \mathbb{F}} x_{ij} \geq r_j \quad \forall j \in \mathbb{C} \\
& && x_{ij} \geq 0, y_i \geq 0 \quad \forall i \in \mathbb{F}, j \in \mathbb{C}
\end{aligned} \tag{1}$$

The dual program is:

$$\begin{aligned}
& \text{maximize} && \sum_{j \in \mathbb{C}} r_j \alpha_j \\
& \text{subject to} && \sum_{j \in \mathbb{C}} \beta_{ij} \leq f_i \quad \forall i \in \mathbb{F} \\
& && \alpha_j - \beta_{ij} \leq d_{ij} \quad \forall i \in \mathbb{F}, j \in \mathbb{C} \\
& && \alpha_j \geq 0, \beta_{ij} \geq 0 \quad \forall i \in \mathbb{F}, j \in \mathbb{C}
\end{aligned} \tag{2}$$

In each of our algorithms we will fix some optimal solutions of the LPs (1) and (2) that we will denote by  $(\mathbf{x}^*, \mathbf{y}^*)$  and  $(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*)$ , respectively.

With  $(\mathbf{x}^*, \mathbf{y}^*)$  fixed, we can define the optimal facility cost as  $F^* = \sum_{i \in \mathbb{F}} f_i y_i^*$  and the optimal connection cost as  $C^* = \sum_{i \in \mathbb{F}, j \in \mathbb{C}} d_{ij} x_{ij}^*$ . Then  $\text{LP}^* = \text{cost}(\mathbf{x}^*, \mathbf{y}^*) = F^* + C^*$  is the joint optimal value of (1) and (2). We can also associate with each client  $j$  its fractional connection cost  $C_j^* = \sum_{i \in \mathbb{F}} d_{ij} x_{ij}^*$ . Clearly,  $C^* = \sum_{j \in \mathbb{C}} C_j^*$ . Throughout the paper we will use notation  $\text{OPT}$  for the optimal integral solution of (1).  $\text{OPT}$  is the value we wish to approximate, but, since  $\text{OPT} \geq \text{LP}^*$ , we can instead use  $\text{LP}^*$  to estimate the approximation ratio of our algorithms.

**Completeness and facility splitting.** Define  $(\mathbf{x}^*, \mathbf{y}^*)$  to be *complete* if  $x_{ij}^* > 0$  implies that  $x_{ij}^* = y_i^*$  for all  $i, j$ . In other words, each connection either uses a site fully or not at all. As shown by Chudak and Shmoys [5], we can modify the given instance by adding at most  $|\mathbb{C}|$  sites to obtain an equivalent instance that has a complete optimal solution, where “equivalent” means that the values of  $F^*$ ,  $C^*$  and  $\text{LP}^*$  are not affected. Roughly, the argument is this: We notice that, without loss of generality, for each client  $k$  there exists at most one site  $i$  such that  $0 < x_{ik}^* < y_i^*$ . We can then perform the following *facility splitting* operation on  $i$ : introduce a new site  $i'$ , let  $y_{i'}^* = y_i^* - x_{ik}^*$ , redefine  $y_i^*$  to be  $x_{ik}^*$ , and then for each client  $j$  redistribute  $x_{ij}^*$  so that  $i$  retains as much connection value as possible and  $i'$  receives the rest. Specifically, we set

$$\begin{aligned}
y_{i'}^* &\leftarrow y_i^* - x_{ik}^*, \quad y_i^* \leftarrow x_{ik}^*, \quad \text{and} \\
x_{i'j}^* &\leftarrow \max(x_{ij}^* - x_{ik}^*, 0), \quad x_{ij}^* \leftarrow \min(x_{ij}^*, x_{ik}^*) \quad \text{for all } j \neq k.
\end{aligned}$$

This operation eliminates the partial connection between  $k$  and  $i$  and does not create any new partial connections. Each client can split at most one site and hence we shall have at most  $|\mathbb{C}|$  more sites.

By the above paragraph, without loss of generality we can assume that the optimal fractional solution  $(\mathbf{x}^*, \mathbf{y}^*)$  is complete. This assumption will in fact greatly simplify some of the arguments in the paper. Additionally, we will frequently use the facility splitting operation described above in our algorithms to obtain fractional solutions with desirable properties.

### 3 Reduction to Polynomial Demands

This section presents a *demand reduction* trick that reduces the problem for arbitrary demands to a special case where demands are bounded by  $|\mathbb{F}|$ , the number of sites. (The formal statement is a little more technical – see Theorem 2.) Our algorithms in the sections that follow process individual demands of each client one by one, and thus they critically rely on the demands being bounded polynomially in terms of  $|\mathbb{F}|$  and  $|\mathbb{C}|$  to keep the overall running time polynomial.

The reduction is based on an optimal fractional solution  $(\mathbf{x}^*, \mathbf{y}^*)$  of LP (1). From the optimality of this solution, we can also assume that  $\sum_{i \in \mathbb{F}} x_{ij}^* = r_j$  for all  $j \in \mathbb{C}$ . As explained in Section 2, we can assume that  $(\mathbf{x}^*, \mathbf{y}^*)$  is complete, that is  $x_{ij}^* > 0$  implies  $x_{ij}^* = y_i^*$  for all  $i, j$ . We split this solution into two parts, namely  $(\mathbf{x}^*, \mathbf{y}^*) = (\hat{\mathbf{x}}, \hat{\mathbf{y}}) + (\dot{\mathbf{x}}, \dot{\mathbf{y}})$ , where

$$\begin{aligned} \hat{y}_i &\leftarrow \lfloor y_i^* \rfloor, & \hat{x}_{ij} &\leftarrow \lfloor x_{ij}^* \rfloor & \text{ and} \\ \dot{y}_i &\leftarrow y_i^* - \lfloor y_i^* \rfloor, & \dot{x}_{ij} &\leftarrow x_{ij}^* - \lfloor x_{ij}^* \rfloor \end{aligned}$$

for all  $i, j$ . Now we construct two FTFP instances  $\hat{\mathcal{I}}$  and  $\dot{\mathcal{I}}$  with the same parameters as the original instance, except that the demand of each client  $j$  is  $\hat{r}_j = \sum_{i \in \mathbb{F}} \hat{x}_{ij}$  in instance  $\hat{\mathcal{I}}$  and  $\dot{r}_j = \sum_{i \in \mathbb{F}} \dot{x}_{ij} = r_j - \hat{r}_j$  in instance  $\dot{\mathcal{I}}$ . It is obvious that if we have integral solutions to both  $\hat{\mathcal{I}}$  and  $\dot{\mathcal{I}}$  then, when added together, they form an integral solution to the original instance. Moreover, we have the following lemma.

**Lemma 1.** (i)  $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$  is a feasible integral solution to instance  $\hat{\mathcal{I}}$ .

(ii)  $(\dot{\mathbf{x}}, \dot{\mathbf{y}})$  is a feasible fractional solution to instance  $\dot{\mathcal{I}}$ .

(iii)  $\dot{r}_j \leq |\mathbb{F}|$  for every client  $j$ .

*Proof.* (i) For feasibility, we need to verify that the constraints of LP (1) are satisfied. Directly from the definition, we have  $\hat{r}_j = \sum_{i \in \mathbb{F}} \hat{x}_{ij}$ . For any  $i$  and  $j$ , by the feasibility of  $(\mathbf{x}^*, \mathbf{y}^*)$  we have  $\hat{x}_{ij} = \lfloor x_{ij}^* \rfloor \leq \lfloor y_i^* \rfloor = \hat{y}_i$ .

(ii) From the definition, we have  $\dot{r}_j = \sum_{i \in \mathbb{F}} \dot{x}_{ij}$ . It remains to show that  $\dot{y}_i \geq \dot{x}_{ij}$  for all  $i, j$ . If  $x_{ij}^* = 0$ , then  $\dot{x}_{ij} = 0$  and we are done. Otherwise, by completeness, we have  $x_{ij}^* = y_i^*$ . Then  $\dot{y}_i = y_i^* - \lfloor y_i^* \rfloor = x_{ij}^* - \lfloor x_{ij}^* \rfloor = \dot{x}_{ij}$ .

(iii) From the definition of  $\dot{x}_{ij}$  we have  $\dot{x}_{ij} < 1$ . Then the bound follows from the definition of  $\dot{r}_j$ .  $\square$

Notice that our construction relies on the completeness assumption; in fact, it is easy to give an example where  $(\dot{\mathbf{x}}, \dot{\mathbf{y}})$  would not be feasible if we used a non-complete optimal solution  $(\mathbf{x}^*, \mathbf{y}^*)$ . Note also that the solutions  $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$  and  $(\dot{\mathbf{x}}, \dot{\mathbf{y}})$  are in fact optimal for their corresponding instances, for if a better solution to  $\hat{\mathcal{I}}$  or  $\dot{\mathcal{I}}$  existed, it could give us a solution to  $\mathcal{I}$  with a smaller objective value.

**Theorem 2.** Suppose that there is a polynomial-time algorithm  $\mathcal{A}$  that, for any instance of FTFP with maximum demand bounded by  $|\mathbb{F}|$ , computes an integral solution that approximates the fractional optimum of this instance within factor  $\rho \geq 1$ . Then there is a  $\rho$ -approximation algorithm  $\mathcal{A}'$  for FTFP.

*Proof.* Given an FTFP instance with arbitrary demands, Algorithm  $\mathcal{A}'$  works as follows: it solves the LP (1) to obtain a fractional optimal solution  $(\mathbf{x}^*, \mathbf{y}^*)$ , then it constructs instances  $\hat{\mathcal{I}}$  and

$\hat{\mathcal{I}}$  described above, applies algorithm  $\mathcal{A}$  to  $\hat{\mathcal{I}}$ , and finally combines (by adding the values) the integral solution  $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$  of  $\hat{\mathcal{I}}$  and the integral solution of  $\hat{\mathcal{I}}$  produced by  $\mathcal{A}$ . This clearly produces a feasible integral solution for the original instance  $\mathcal{I}$ . The solution produced by  $\mathcal{A}$  has cost at most  $\rho \cdot \text{cost}(\hat{\mathbf{x}}, \hat{\mathbf{y}})$ , because  $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$  is feasible for  $\hat{\mathcal{I}}$ . Thus the cost of  $\mathcal{A}'$  is at most

$$\text{cost}(\hat{\mathbf{x}}, \hat{\mathbf{y}}) + \rho \cdot \text{cost}(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \leq \rho(\text{cost}(\hat{\mathbf{x}}, \hat{\mathbf{y}}) + \text{cost}(\hat{\mathbf{x}}, \hat{\mathbf{y}})) = \rho \cdot \text{LP}^* \leq \rho \cdot \text{OPT},$$

where the first inequality follows from  $\rho \geq 1$ . This completes the proof.  $\square$

## 4 Adaptive Partitioning

In this section we develop our second technique, which we call *adaptive partitioning*. Given an FTFP instance and an optimal fractional solution  $(\mathbf{x}^*, \mathbf{y}^*)$  to LP (1), we split each client  $j$  into  $r_j$  individual *unit demand points* (or just *demands*), and we split each site  $i$  into no more than  $|\mathbb{F}| + 2R|\mathbb{C}|^2$  *facility points* (or *facilities*), where  $R = \max_{j \in \mathbb{C}} r_j$ . We denote the demand set by  $\overline{\mathbb{C}}$  and the facility set by  $\overline{\mathbb{F}}$ , respectively. We will also partition  $(\mathbf{x}^*, \mathbf{y}^*)$  into a fractional solution  $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$  for the split instance. We will typically use symbols  $\nu$  and  $\mu$  to index demands and facilities respectively, that is  $\bar{\mathbf{x}} = (\bar{x}_{\mu\nu})$  and  $\bar{\mathbf{y}} = (\bar{y}_\mu)$ . As before, the *neighborhood of a demand*  $\nu$  is  $\overline{N}(\nu) = \{\mu \in \overline{\mathbb{F}} : \bar{x}_{\mu\nu} > 0\}$ . We will use notation  $\nu \in j$  to mean that  $\nu$  is a demand of client  $j$ ; similarly,  $\mu \in i$  means that facility  $\mu$  is on site  $i$ . Different demands of the same client (that is,  $\nu, \nu' \in j$ ) are called *siblings*. Further, we use the convention that  $f_\mu = f_i$  for  $\mu \in i$ ,  $\alpha_\nu^* = \alpha_j^*$  for  $\nu \in j$  and  $d_{\mu\nu} = d_{\mu j} = d_{ij}$  for  $\mu \in i$  and  $\nu \in j$ . We define  $C_\nu^{\text{avg}} = \sum_{\mu \in \overline{N}(\nu)} d_{\mu\nu} \bar{x}_{\mu\nu} = \sum_{\mu \in \overline{\mathbb{F}}} d_{\mu\nu} \bar{x}_{\mu\nu}$ . One can think of  $C_\nu^{\text{avg}}$  as the average connection cost of demand  $\nu$ , if we chose a connection to facility  $\mu$  with probability  $\bar{x}_{\mu\nu}$ . In our partitioned fractional solution we guarantee for every  $\nu$  that  $\sum_{\mu \in \overline{\mathbb{F}}} \bar{x}_{\mu\nu} = 1$ .

Some demands in  $\overline{\mathbb{C}}$  will be designated as *primary demands* and the set of primary demands will be denoted by  $P$ . In addition, we will use the overlap structure between demand neighborhoods to define a mapping that assigns each demand  $\nu \in \overline{\mathbb{C}}$  to some primary demand  $\kappa \in P$ . As shown in the rounding algorithms in later sections, for each primary demand we guarantee exactly one open facility in its neighborhood, while for a non-primary demand, there is constant probability that none of its neighbors open. In this case we estimate its connection cost by the distance to the facility opened in its assigned primary demand's neighborhood. For this reason the connection cost of a primary demand must be “small” compared to the non-primary demands assigned to it. We also need sibling demands assigned to different primary demands to satisfy the fault-tolerance requirement. Specifically, this partitioning will be constructed to satisfy a number of properties that are detailed below.

(PS) *Partitioned solution.* Vector  $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$  is a partition of  $(\mathbf{x}^*, \mathbf{y}^*)$ , with unit-value demands, that is:

1.  $\sum_{\mu \in \overline{\mathbb{F}}} \bar{x}_{\mu\nu} = 1$  for each demand  $\nu \in \overline{\mathbb{C}}$ .
2.  $\sum_{\mu \in i, \nu \in j} \bar{x}_{\mu\nu} = x_{ij}^*$  for each site  $i \in \mathbb{F}$  and client  $j \in \mathbb{C}$ .
3.  $\sum_{\mu \in i} \bar{y}_\mu = y_i^*$  for each site  $i \in \mathbb{F}$ .

(CO) *Completeness.* Solution  $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$  is complete, that is  $\bar{x}_{\mu\nu} \neq 0$  implies  $\bar{x}_{\mu\nu} = \bar{y}_\mu$ , for all  $\mu \in \overline{\mathbb{F}}, \nu \in \overline{\mathbb{C}}$ .

(PD) *Primary demands.* Primary demands satisfy the following conditions:

1. For any two different primary demands  $\kappa, \kappa' \in P$  we have  $\overline{N}(\kappa) \cap \overline{N}(\kappa') = \emptyset$ .
2. For each site  $i \in \mathbb{F}$ ,  $\sum_{\mu \in i} \sum_{\kappa \in P} \bar{x}_{\mu\kappa} \leq y_i^*$ .
3. Each demand  $\nu \in \overline{\mathbb{C}}$  is assigned to one primary demand  $\kappa \in P$  such that
  - (a)  $\overline{N}(\nu) \cap \overline{N}(\kappa) \neq \emptyset$ , and
  - (b)  $C_\nu^{\text{avg}} + \alpha_\nu^* \geq C_\kappa^{\text{avg}} + \alpha_\kappa^*$ .

(SI) *Siblings.* For any pair  $\nu, \nu'$  of different siblings we have

1.  $\overline{N}(\nu) \cap \overline{N}(\nu') = \emptyset$ .
2. If  $\nu$  is assigned to a primary demand  $\kappa$  then  $\overline{N}(\nu') \cap \overline{N}(\kappa) = \emptyset$ . In particular, by Property PD(3(a)), this implies that different sibling demands are assigned to different primary demands.

As we shall demonstrate in later sections, these properties allow us to extend known UFL rounding algorithms to obtain an integral solution to our FTFP problem with a matching approximation ratio. Our partitioning is “adaptive” in the sense that it is constructed one demand at a time, and the connection values for the demands of a client depend on the choice of earlier demands, of this or other clients, and their connection values. We would like to point out that the adaptive partitioning process for the 1.575-approximation algorithm (Section 7) is more subtle than the 3-approximation (Section 5) and the 1.736-approximation algorithms (Section 6), due to the introduction of close and far neighborhood.

**Implementation of Adaptive Partitioning.** We now describe an algorithm for partitioning the instance and the fractional solution so that the properties (PS), (CO), (PD), and (SI) are satisfied. Recall that  $\overline{\mathbb{F}}$  and  $\overline{\mathbb{C}}$ , respectively, denote the sets of facilities and demands that will be created in this stage, and  $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$  is the partitioned solution to be computed.

The adaptive partitioning algorithm consists of two phases: Phase 1 is called the partition phase and Phase 2 is called the augmenting phase. Phase 1 is done in iterations, where in each iteration we find the “best” client  $j$  and create a new demand  $\nu$  out of it. This demand either becomes a primary demand itself, or it is assigned to some existing primary demand. We call a client  $j$  *exhausted* when all its  $r_j$  demands have been created and assigned to some primary demands. Phase 1 completes when all clients are exhausted. In Phase 2 we ensure that every demand has a total connection values equal to 1, that is condition (PS.1).

For each site  $i$  we will initially create one “big” facility  $\mu$  with initial value  $\bar{y}_\mu = y_i^*$ . While we partition the instance, creating new demands and connections, this facility may end up being split into more facilities to preserve completeness of the fractional solution. Also, we will gradually decrease the fractional connection vector for each client  $j$ , to account for the demands already created for  $j$  and their connection values. These decreased connection values will be stored in an auxiliary vector  $\tilde{\mathbf{x}}$ . The intuition is that  $\tilde{\mathbf{x}}$  represents the part of  $\mathbf{x}^*$  that still has not been allocated to existing demands and future demands can use  $\tilde{\mathbf{x}}$  for their connections. For technical reasons,  $\tilde{\mathbf{x}}$  will be indexed by facilities (rather than sites) and clients, that is  $\tilde{\mathbf{x}} = (\tilde{x}_{\mu j})$ . At the beginning, we set  $\tilde{x}_{\mu j} \leftarrow x_{ij}^*$  for each  $j \in \mathbb{C}$ , where  $\mu \in i$  is the single facility created initially at site  $i$ . At each step, whenever we create a new demand  $\nu$  for a client  $j$ , we will define its values  $\bar{x}_{\mu\nu}$  and appropriately

reduce the values  $\tilde{x}_{\mu j}$ , for all facilities  $\mu$ . We will deal with two types of neighborhoods, with respect to  $\tilde{\mathbf{x}}$  and  $\bar{\mathbf{x}}$ , that is  $\tilde{N}(j) = \{\mu \in \bar{\mathbb{F}} : \tilde{x}_{\mu j} > 0\}$  for  $j \in \mathbb{C}$  and  $\bar{N}(\nu) = \{\mu \in \bar{\mathbb{F}} : \bar{x}_{\mu\nu} > 0\}$  for  $\nu \in \bar{\mathbb{C}}$ . During this process we preserve the completeness (CO) of the fractional solutions  $\tilde{\mathbf{x}}$  and  $\bar{\mathbf{x}}$ . More precisely, the following properties will hold for every facility  $\mu$  after every iteration:

- (c1) For each demand  $\nu$  either  $\bar{x}_{\mu\nu} = 0$  or  $\bar{x}_{\mu\nu} = \bar{y}_\mu$ . This is the same condition as condition (CO), yet we repeat it here as (c1) needs to hold after every iteration, while condition (CO) only applies to the final partitioned fractional solution  $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$ .
- (c2) For each client  $j$ , either  $\tilde{x}_{\mu j} = 0$  or  $\tilde{x}_{\mu j} = \bar{y}_\mu$ .

A full description of the algorithm is given in Pseudocode 1. Initially, the set  $U$  of non-exhausted clients contains all clients, the set  $\bar{\mathbb{C}}$  of demands is empty, the set  $\bar{\mathbb{F}}$  of facilities consists of one facility  $\mu$  on each site  $i$  with  $\bar{y}_\mu = y_i^*$ , and the set  $P$  of primary demands is empty (Lines 1–4). In one iteration of the while loop (Lines 5–8), for each client  $j$  we compute a quantity called  $\text{tcc}(j)$  (tentative connection cost), that represents the average distance from  $j$  to the set  $\tilde{N}_1(j)$  of the nearest facilities  $\mu$  whose total connection value to  $j$  (the sum of  $\tilde{x}_{\mu j}$ 's) equals 1. This set is computed by Procedure NEARESTUNITCHUNK() (see Pseudocode 2, Lines 1–9), which adds facilities to  $\tilde{N}_1(j)$  in order of nondecreasing distance, until the total connection value is exactly 1. (The procedure actually uses the  $\bar{y}_\mu$  values, which are equal to the connection values, by the completeness condition (c2).) This may require splitting the last added facility and adjusting the connection values so that conditions (c1) and (c2) are preserved.

The next step is to pick a client  $p$  with minimum  $\text{tcc}(p) + \alpha_p^*$  and create a demand  $\nu$  for  $p$  (Lines 9–10). If  $\tilde{N}_1(p)$  overlaps the neighborhood of some existing primary demand  $\kappa$  (if there are multiple such  $\kappa$ 's, pick any of them), we assign  $\nu$  to  $\kappa$ , and  $\nu$  acquires all the connection values  $\tilde{x}_{\mu p}$  between client  $p$  and facility  $\mu$  in  $\tilde{N}(p) \cap \bar{N}(\kappa)$  (Lines 11–13). Note that although we check for overlap with  $\tilde{N}_1(p)$ , we then move all facilities in the intersection with  $\tilde{N}(p)$ , a bigger set, into  $\bar{N}(\nu)$ . The other case is when  $\tilde{N}_1(p)$  is disjoint from the neighborhoods of all existing primary demands. Then, in Lines 15–16,  $\nu$  becomes itself a primary demand and we assign  $\nu$  to itself. It also inherits the connection values to all facilities  $\mu \in \tilde{N}_1(p)$  from  $p$  (recall that  $\tilde{x}_{\mu p} = \bar{y}_\mu$ ), with all other  $\bar{x}_{\mu\nu}$  values set to 0.

At this point all primary demands satisfy Property (PS.1), but this may not be true for non-primary demands. For those demands we still may need to adjust the  $\bar{x}_{\mu\nu}$  values so that the total connection value for  $\nu$ , that is  $\text{conn}(\nu) \stackrel{\text{def}}{=} \sum_{\mu \in \bar{\mathbb{F}}} \bar{x}_{\mu\nu}$ , is equal 1. This is accomplished by Procedure AUGMENTTOUNIT() (definition in Pseudocode 2, Lines 10–21) that allocates to  $\nu \in j$  some of the remaining connection values  $\tilde{x}_{\mu j}$  of client  $j$  (Lines 19–21). AUGMENTTOUNIT() will repeatedly pick any facility  $\eta$  with  $\tilde{x}_{\eta j} > 0$ . If  $\tilde{x}_{\eta j} \leq 1 - \text{conn}(\nu)$ , then the connection value  $\tilde{x}_{\eta j}$  is reassigned to  $\nu$ . Otherwise,  $\tilde{x}_{\eta j} > 1 - \text{conn}(\nu)$ , in which case we split  $\eta$  so that connecting  $\nu$  to one of the created copies of  $\eta$  will make  $\text{conn}(\nu)$  equal 1, and we'll be done.

Notice that we start with  $|\mathbb{F}|$  facilities and in each iteration of the while loop in Line 5 each client causes at most one split. We have a total of no more than  $R|\mathbb{C}|$  iterations as in each iteration we create one demand. (Recall that  $R = \max_j r_j$ .) In Phase 2 we do an augment step for each demand  $\nu$  and this creates no more than  $R|\mathbb{C}|$  new facilities. So the total number of facilities we created will be at most  $|\mathbb{F}| + R|\mathbb{C}|^2 + R|\mathbb{C}| \leq |\mathbb{F}| + 2R|\mathbb{C}|^2$ , which is polynomial in  $|\mathbb{F}| + |\mathbb{C}|$  due to our earlier bound on  $R$ .



---

**Pseudocode 1** Algorithm: Adaptive Partitioning
 

---

**Input:**  $\mathbb{F}, \mathbb{C}, (\mathbf{x}^*, \mathbf{y}^*)$ 
**Output:**  $\bar{\mathbb{F}}, \bar{\mathbb{C}}, (\bar{\mathbf{x}}, \bar{\mathbf{y}})$ 
 $\triangleright$  Unspecified  $\bar{x}_{\mu\nu}$ 's and  $\tilde{x}_{\mu j}$ 's are assumed to be 0

```

1:  $\tilde{r} \leftarrow \mathbf{r}, U \leftarrow \mathbb{C}, \bar{\mathbb{F}} \leftarrow \emptyset, \bar{\mathbb{C}} \leftarrow \emptyset, P \leftarrow \emptyset$   $\triangleright$  Phase 1
2: for each site  $i \in \mathbb{F}$  do
3:   create a facility  $\mu$  at  $i$  and add  $\mu$  to  $\bar{\mathbb{F}}$ 
4:    $\bar{y}_\mu \leftarrow y_i^*$  and  $\tilde{x}_{\mu j} \leftarrow x_{ij}^*$  for each  $j \in \mathbb{C}$ 
5: while  $U \neq \emptyset$  do
6:   for each  $j \in U$  do
7:      $\tilde{N}_1(j) \leftarrow \text{NEARESTUNITCHUNK}(j, \bar{\mathbb{F}}, \tilde{\mathbf{x}}, \bar{\mathbf{x}}, \bar{\mathbf{y}})$   $\triangleright$  see Pseudocode 2
8:      $\text{tcc}(j) \leftarrow \sum_{\mu \in \tilde{N}_1(j)} d_{\mu j} \cdot \tilde{x}_{\mu j}$ 
9:      $p \leftarrow \arg \min_{j \in U} \{\text{tcc}(j) + \alpha_j^*\}$ 
10:    create a new demand  $\nu$  for client  $p$ 
11:    if  $\tilde{N}_1(p) \cap \bar{N}(\kappa) \neq \emptyset$  for some primary demand  $\kappa \in P$  then
12:      assign  $\nu$  to  $\kappa$ 
13:       $\bar{x}_{\mu\nu} \leftarrow \tilde{x}_{\mu p}$  and  $\tilde{x}_{\mu p} \leftarrow 0$  for each  $\mu \in \tilde{N}(p) \cap \bar{N}(\kappa)$ 
14:    else
15:      make  $\nu$  primary,  $P \leftarrow P \cup \{\nu\}$ , assign  $\nu$  to itself
16:      set  $\bar{x}_{\mu\nu} \leftarrow \tilde{x}_{\mu p}$  and  $\tilde{x}_{\mu p} \leftarrow 0$  for each  $\mu \in \tilde{N}_1(p)$ 
17:       $\bar{\mathbb{C}} \leftarrow \bar{\mathbb{C}} \cup \{\nu\}, \bar{r}_p \leftarrow \bar{r}_p - 1$ 
18:      if  $\bar{r}_p = 0$  then  $U \leftarrow U \setminus \{p\}$ 
19: for each client  $j \in \mathbb{C}$  do  $\triangleright$  Phase 2
20:   for each demand  $\nu \in j$  do  $\triangleright$  each client  $j$  has  $r_j$  demands
21:    if  $\sum_{\mu \in \bar{N}(\nu)} \bar{x}_{\mu\nu} < 1$  then  $\text{AUGMENTTOUNIT}(\nu, j, \bar{\mathbb{F}}, \tilde{\mathbf{x}}, \bar{\mathbf{x}}, \bar{\mathbf{y}})$   $\triangleright$  see Pseudocode 2

```

---

*Correctness.* We now show that all the required properties (PS), (CO), (PD) and (SI) are satisfied by the above construction.

Properties (PS) and (CO) follow directly from the algorithm. (CO) is implied by the completeness condition (c1) that the algorithm maintains after each iteration. Condition (PS.1) is a result of calling Procedure AUGMENTTOUNIT() in Line 21. To see that (PS.2) holds, note that at each step the algorithm maintains the invariant that, for every  $i \in \mathbb{F}$  and  $j \in \mathbb{C}$ , we have  $\sum_{\mu \in i} \sum_{\nu \in j} \bar{x}_{\mu\nu} + \sum_{\mu \in i} \tilde{x}_{\mu j} = x_{ij}^*$ . In the end, we will create  $r_j$  demands for each client  $j$ , with each demand  $\nu \in j$  satisfying (PS.1), and thus  $\sum_{\nu \in j} \sum_{\mu \in \bar{\mathbb{F}}} \bar{x}_{\mu\nu} = r_j$ . This implies that  $\tilde{x}_{\mu j} = 0$  for every facility  $\mu \in \bar{\mathbb{F}}$ , and PS(2) follows. PS(3) holds because every time we split a facility  $\mu$  into  $\mu'$  and  $\mu''$ , the sum of  $\bar{y}_{\mu'}$  and  $\bar{y}_{\mu''}$  is equal to the old value of  $\bar{y}_\mu$ .

Now we deal with properties in group (PD). First, (PD.1) follows directly from the algorithm, Pseudocode 1 (Lines 14–16), since every primary demand has its neighborhood fixed when created, and that neighborhood is disjoint from those of the existing primary demands.

Property (PD.2) follows from (PD.1), (CO) and (PS.3). In more detail, it can be justified as follows. By (PD.1), for each  $\mu \in i$  there is at most one  $\kappa \in P$  with  $\bar{x}_{\mu\kappa} > 0$  and we have  $\bar{x}_{\mu\kappa} = \bar{y}_\mu$  due to (CO). Let  $K \subseteq i$  be the set of those  $\mu$ 's for which such  $\kappa \in P$  exists, and denote this  $\kappa$  by  $\kappa_\mu$ . Then, using conditions (CO) and (PS.3), we have  $\sum_{\mu \in i} \sum_{\kappa \in P} \bar{x}_{\mu\kappa} = \sum_{\mu \in K} \bar{x}_{\mu\kappa_\mu} = \sum_{\mu \in K} \bar{y}_\mu \leq \sum_{\mu \in i} \bar{y}_\mu = y_i^*$ .

---

**Pseudocode 2** Helper functions used in Pseudocode 1
 

---

```

1: function NEARESTUNITCHUNK( $j, \bar{\mathbb{F}}, \tilde{\mathbf{x}}, \bar{\mathbf{x}}, \bar{\mathbf{y}}$ )                                ▷ upon return,  $\sum_{\mu \in \tilde{N}_1(j)} \tilde{x}_{\mu j} = 1$ 
2:   Let  $\tilde{N}(j) = \{\mu_1, \dots, \mu_q\}$  where  $d_{\mu_1 j} \leq d_{\mu_2 j} \leq \dots \leq d_{\mu_q j}$ 
3:   Let  $l$  be such that  $\sum_{k=1}^l \bar{y}_{\mu_k} \geq 1$  and  $\sum_{k=1}^{l-1} \bar{y}_{\mu_k} < 1$ 
4:   Create a new facility  $\sigma$  at the same site as  $\mu_l$  and add it to  $\bar{\mathbb{F}}$                                 ▷ split  $\mu_l$ 
5:   Set  $\bar{y}_\sigma \leftarrow \sum_{k=1}^l \bar{y}_{\mu_k} - 1$  and  $\bar{y}_{\mu_l} \leftarrow \bar{y}_{\mu_l} - \bar{y}_\sigma$ 
6:   For each  $\nu \in \mathbb{C}$  with  $\bar{x}_{\mu_l \nu} > 0$  set  $\bar{x}_{\mu_l \nu} \leftarrow \bar{y}_{\mu_l}$  and  $\bar{x}_{\sigma \nu} \leftarrow \bar{y}_\sigma$ 
7:   For each  $j' \in \mathbb{C}$  with  $\tilde{x}_{\mu_l j'} > 0$  (including  $j$ ) set  $\tilde{x}_{\mu_l j'} \leftarrow \bar{y}_{\mu_l}$  and  $\tilde{x}_{\sigma j'} \leftarrow \bar{y}_\sigma$ 
8:   (All other new connection values are set to 0)
9:   return  $\tilde{N}_1(j) = \{\mu_1, \dots, \mu_{l-1}, \mu_l\}$ 

10: function AUGMENTTOUNIT( $\nu, j, \bar{\mathbb{F}}, \tilde{\mathbf{x}}, \bar{\mathbf{x}}, \bar{\mathbf{y}}$ )                                ▷  $\nu$  is a demand of client  $j$ 
11:   while  $\sum_{\mu \in \bar{\mathbb{F}}} \bar{x}_{\mu \nu} < 1$  do                                ▷ upon return,  $\sum_{\mu \in \bar{N}(\nu)} \bar{x}_{\mu \nu} = 1$ 
12:     Let  $\eta$  be any facility such that  $\tilde{x}_{\eta j} > 0$ 
13:     if  $1 - \sum_{\mu \in \bar{\mathbb{F}}} \bar{x}_{\mu \nu} \geq \tilde{x}_{\eta j}$  then
14:        $\bar{x}_{\eta \nu} \leftarrow \tilde{x}_{\eta j}, \tilde{x}_{\eta j} \leftarrow 0$ 
15:     else
16:       Create a new facility  $\sigma$  at the same site as  $\eta$  and add it to  $\bar{\mathbb{F}}$                                 ▷ split  $\eta$ 
17:       Let  $\bar{y}_\sigma \leftarrow 1 - \sum_{\mu \in \bar{\mathbb{F}}} \bar{x}_{\mu \nu}, \bar{y}_\eta \leftarrow \bar{y}_\eta - \bar{y}_\sigma$ 
18:       Set  $\bar{x}_{\sigma \nu} \leftarrow \bar{y}_\sigma, \bar{x}_{\eta \nu} \leftarrow 0, \tilde{x}_{\eta j} \leftarrow \bar{y}_\eta, \tilde{x}_{\sigma j} \leftarrow 0$ 
19:       For each  $\nu' \neq \nu$  with  $\bar{x}_{\eta \nu'} > 0$  set  $\bar{x}_{\eta \nu'} \leftarrow \bar{y}_\eta, \bar{x}_{\sigma \nu'} \leftarrow \bar{y}_\sigma$ 
20:       For each  $j' \neq j$  with  $\tilde{x}_{\eta j'} > 0$  set  $\tilde{x}_{\eta j'} \leftarrow \bar{y}_\eta, \tilde{x}_{\sigma j'} \leftarrow \bar{y}_\sigma$ 
21:       (All other new connection values are set to 0)
  
```

---

Property (PD.3(a)) follows from the way the algorithm assigns primary demands. When demand  $\nu$  of client  $p$  is assigned to a primary demand  $\kappa$  in Lines 11–13 of Pseudocode 1, we move all facilities in  $\tilde{N}(p) \cap \bar{N}(\kappa)$  (the intersection is nonempty) into  $\bar{N}(\nu)$ , and we never remove a facility from  $\bar{N}(\nu)$ . We postpone the proof for (PD.3(b)) to Lemma 5.

Finally we argue that the properties in group (SI) hold. (SI.1) is easy, since for any client  $j$ , each facility  $\mu$  is added to the neighborhood of at most one demand  $\nu \in j$ , by setting  $\bar{x}_{\mu \nu}$  to  $\bar{y}_\mu$ , while other siblings  $\nu'$  of  $\nu$  have  $\bar{x}_{\mu \nu'} = 0$ . Note that right after a demand  $\nu \in p$  is created, its neighborhood is disjoint from the neighborhood of  $p$ , that is  $\bar{N}(\nu) \cap \tilde{N}(p) = \emptyset$ , by Lines 11–13 of the algorithm. Thus all demands of  $p$  created later will have neighborhoods disjoint from the set  $\bar{N}(\nu)$  before the augmenting phase 2. Furthermore, Procedure AUGMENTTOUNIT() preserves this property, because when it adds a facility to  $\bar{N}(\nu)$  then it removes it from  $\tilde{N}(p)$ , and in case of splitting, one resulting facility is added to  $\bar{N}(\nu)$  and the other to  $\tilde{N}(p)$ . Property (SI.2) is shown below in Lemma 3.

It remains to show Properties (PD.3(b)) and (SI.2). We show them in the lemmas below, thus completing the description of our adaptive partition process.

**Lemma 3.** *Property (SI.2) holds after the Adaptive Partitioning stage.*

*Proof.* Let  $\nu_1, \dots, \nu_{r_j}$  be the demands of a client  $j \in \mathbb{C}$ , listed in the order of creation, and, for each  $q = 1, 2, \dots, r_j$ , denote by  $\kappa_q$  the primary demand that  $\nu_q$  is assigned to. After the completion of Phase 1 of Pseudocode 1 (Lines 5–18), we have  $\bar{N}(\nu_s) \subseteq \bar{N}(\kappa_s)$  for  $s = 1, \dots, r_j$ . Since any two

primary demands have disjoint neighborhoods, we have  $\overline{N}(\nu_s) \cap \overline{N}(\kappa_q) = \emptyset$  for any  $s \neq q$ , that is Property (SI.2) holds right after Phase 1.

After Phase 1 all neighborhoods  $\overline{N}(\kappa_s)$ ,  $s = 1, \dots, r_j$  have already been fixed and they do not change in Phase 2. None of the facilities in  $\tilde{N}(j)$  appear in any of  $\overline{N}(\kappa_s)$  for  $s = 1, \dots, r_j$ , by the way we allocate facilities in Lines 13 and 16. Therefore during the augmentation process in Phase 2, when we add facilities from  $\tilde{N}(j)$  to  $\overline{N}(\nu)$ , for some  $\nu \in j$  (Line 19–21 of Pseudocode 1), all the required disjointness conditions will be preserved.  $\square$

We need one more lemma before proving our last property (PD.3(b)). For a client  $j$  and a demand  $\nu$ , we use notation  $\text{tcc}_\nu(j)$  for the value of  $\text{tcc}(j)$  at the time when  $\nu$  was created. (It is not necessary that  $\nu \in j$  but we assume that  $j$  is not exhausted at that time.)

**Lemma 4.** *Let  $\eta$  and  $\nu$  be two demands, with  $\eta$  created not later than  $\nu$ , and let  $j \in \mathbb{C}$  be a client that is not exhausted when  $\nu$  is created. Then we have*

- (a)  $\text{tcc}_\eta(j) \leq \text{tcc}_\nu(j)$ , and
- (b) if  $\nu \in j$  then  $\text{tcc}_\eta(j) \leq C_\nu^{\text{avg}}$ .

*Proof.* We focus first on the time when demand  $\eta$  is about to be created, right after the call to `NEARESTUNITCHUNK()` in Pseudocode 1, Line 7. Let  $\tilde{N}(j) = \{\mu_1, \dots, \mu_q\}$  with all facilities  $\mu_s$  ordered according to nondecreasing distance from  $j$ . Consider the following linear program:

$$\begin{aligned} & \text{minimize} && \sum_s d_{\mu_s j} z_s \\ & \text{subject to} && \sum_s z_s \geq 1 \\ & && 0 \leq z_s \leq \tilde{x}_{\mu_s j} \quad \text{for all } s \end{aligned}$$

This is a fractional minimum knapsack covering problem (with knapsack size equal 1) and its optimal fractional solution is the greedy solution, whose value is exactly  $\text{tcc}_\eta(j)$ .

On the other hand, we claim that  $\text{tcc}_\nu(j)$  can be thought of as the value of some feasible solution to this linear program, and that the same is true for  $C_\nu^{\text{avg}}$  if  $\nu \in j$ . Indeed, each of these quantities involves some later values  $\tilde{x}_{\mu j}$ , where  $\mu$  could be one of the facilities  $\mu_s$  or a new facility obtained from splitting. For each  $s$ , however, the sum of all values  $\tilde{x}_{\mu j}$ , over the facilities  $\mu$  that were split from  $\mu_s$ , cannot exceed the value  $\tilde{x}_{\mu_s j}$  at the time when  $\eta$  was created, because splitting facilities preserves this sum and creating new demands for  $j$  can only decrease it. Therefore both quantities  $\text{tcc}_\nu(j)$  and  $C_\nu^{\text{avg}}$  (for  $\nu \in j$ ) correspond to some choice of the  $z_s$  variables (adding up to 1), and the lemma follows.  $\square$

**Lemma 5.** *Property (PD.3(b)) holds after the Adaptive Partitioning stage.*

*Proof.* Suppose that demand  $\nu \in j$  is assigned to some primary demand  $\kappa \in p$ . Then

$$C_\kappa^{\text{avg}} + \alpha_\kappa^* = \text{tcc}_\kappa(p) + \alpha_p^* \leq \text{tcc}_\kappa(j) + \alpha_j^* \leq C_\nu^{\text{avg}} + \alpha_\nu^*.$$

We now justify this derivation. By definition we have  $\alpha_\kappa^* = \alpha_p^*$ . Further, by the algorithm, if  $\kappa$  is a primary demand of client  $p$ , then  $C_\kappa^{\text{avg}}$  is equal to  $\text{tcc}(p)$  computed when  $\kappa$  is created, which is exactly  $\text{tcc}_\kappa(p)$ . Thus the first equation is true. The first inequality follows from the choice of  $p$  in Line 9 in Pseudocode 1. The last inequality holds because  $\alpha_j^* = \alpha_\nu^*$  (due to  $\nu \in j$ ), and because  $\text{tcc}_\kappa(j) \leq C_\nu^{\text{avg}}$ , which follows from Lemma 4.  $\square$

We have thus proved that all properties (PS), (CO), (PD) and (SI) hold for our partitioned fractional solution  $(\bar{x}, \bar{y})$ . In the following sections we show how to use these properties to round the fractional solution to an approximate integral solution. For the 3-approximation algorithm (Section 5) and the 1.736-approximation algorithm (Section 6), the first phase of the algorithm is exactly the same partition process as described above. However, the 1.575-approximation algorithm (Section 7) demands a more sophisticated partitioning process as the interplay between close and far neighborhood of sibling demands result in more delicate properties that our partitioned fractional solution must satisfy.

## 5 Algorithm EGUP with Ratio 3

With the partitioned FTFP instance and its associated fractional solution in place, we now begin to introduce our rounding algorithms. The algorithm we describe in this section achieves ratio 3. Although this is still quite far from our best ratio 1.575 that we derive later, we include this algorithm in the paper to illustrate, in a relatively simple setting, how the properties of our partitioned fractional solution are used in rounding it to an integral solution with cost not too far away from an optimal solution. The rounding approach we use here is an extension of the corresponding method for UFL described in [8].

**Algorithm EGUP.** At a high level, we would open exactly one facility for each primary demand  $\kappa$ , and each non-primary demand is connected to the facility opened for the primary demand it was assigned to.

More precisely, we apply a rounding process, guided by the fractional values  $(\bar{y}_\mu)$  and  $(\bar{x}_{\mu\nu})$ , that produces an integral solution. This integral solution is obtained by choosing a subset of facilities in  $\bar{\mathbb{F}}$  to open, and for each demand in  $\bar{\mathbb{C}}$ , specifying an open facility that this demand will be connected to. For each primary demand  $\kappa \in P$ , we want to open one facility  $\phi(\kappa) \in \bar{N}(\kappa)$ . To this end, we use randomization: for each  $\mu \in \bar{N}(\kappa)$ , we choose  $\phi(\kappa) = \mu$  with probability  $\bar{x}_{\mu\kappa}$ , ensuring that exactly one  $\mu \in \bar{N}(\kappa)$  is chosen. Note that  $\sum_{\mu \in \bar{N}(\kappa)} \bar{x}_{\mu\kappa} = 1$ , so this distribution is well-defined. We open this facility  $\phi(\kappa)$  and connect to  $\phi(\kappa)$  all demands that are assigned to  $\kappa$ .

In our description above, the algorithm is presented as a randomized algorithm. It can be derandomized using the method of conditional expectations, which is commonly used in approximation algorithms for facility location problems and standard enough that presenting it here would be redundant. Readers less familiar with this field are recommended to consult [5], where the method of conditional expectations is applied in a context very similar to ours.

**Analysis.** We now bound the expected facility cost and connection cost by establishing the two lemmas below.

**Lemma 6.** *The expectation of facility cost  $F_{\text{EGUP}}$  of our solution is at most  $F^*$ .*

*Proof.* By Property (PD.1), the neighborhoods of primary demands are disjoint. Also, for any primary demand  $\kappa \in P$ , the probability that a facility  $\mu \in \bar{N}(\kappa)$  is chosen as the open facility  $\phi(\kappa)$

is  $\bar{x}_{\mu\kappa}$ . Hence the expected total facility cost is

$$\begin{aligned}\text{Exp}[F_{\text{EGUP}}] &= \sum_{\kappa \in P} \sum_{\mu \in \bar{N}(\kappa)} f_{\mu} \bar{x}_{\mu\kappa} \\ &= \sum_{\kappa \in P} \sum_{\mu \in \bar{\mathbb{F}}} f_{\mu} \bar{x}_{\mu\kappa} \\ &= \sum_{i \in \bar{\mathbb{F}}} f_i \sum_{\mu \in i} \sum_{\kappa \in P} \bar{x}_{\mu\kappa} \\ &\leq \sum_{i \in \bar{\mathbb{F}}} f_i y_i^* = F^*,\end{aligned}$$

where the inequality follows from Property (PD.2).  $\square$

**Lemma 7.** *The expectation of connection cost  $C_{\text{EGUP}}$  of our solution is at most  $C^* + 2 \cdot \text{LP}^*$ .*

*Proof.* For a primary demand  $\kappa$ , its expected connection cost is  $C_{\kappa}^{\text{avg}}$  because we choose facility  $\mu$  with probability  $\bar{x}_{\mu\kappa}$ .

Consider a non-primary demand  $\nu$  assigned to a primary demand  $\kappa \in P$ . Let  $\mu$  be any facility in  $\bar{N}(\nu) \cap \bar{N}(\kappa)$ . Since  $\mu$  is in both  $\bar{N}(\nu)$  and  $\bar{N}(\kappa)$ , we have  $d_{\mu\nu} \leq \alpha_{\nu}^*$  and  $d_{\mu\kappa} \leq \alpha_{\kappa}^*$  (This follows from the complementary slackness conditions since  $\alpha_{\nu}^* = \beta_{\mu\nu}^* + d_{\mu\nu}$  for each  $\mu \in \bar{N}(\nu)$ ). Thus, applying the triangle inequality, for any fixed choice of facility  $\phi(\kappa)$  we have

$$d_{\phi(\kappa)\nu} \leq d_{\phi(\kappa)\kappa} + d_{\mu\kappa} + d_{\mu\nu} \leq d_{\phi(\kappa)\kappa} + \alpha_{\kappa}^* + \alpha_{\nu}^*.$$

Therefore the expected distance from  $\nu$  to its facility  $\phi(\kappa)$  is

$$\begin{aligned}\text{Exp}[d_{\phi(\kappa)\nu}] &\leq C_{\kappa}^{\text{avg}} + \alpha_{\kappa}^* + \alpha_{\nu}^* \\ &\leq C_{\nu}^{\text{avg}} + \alpha_{\nu}^* + \alpha_{\nu}^* = C_{\nu}^{\text{avg}} + 2\alpha_{\nu}^*,\end{aligned}$$

where the second inequality follows from Property (PD.3(b)). From the definition of  $C_{\nu}^{\text{avg}}$  and Property (PS.2), for any  $j \in \mathbb{C}$  we have

$$\begin{aligned}\sum_{\nu \in j} C_{\nu}^{\text{avg}} &= \sum_{\nu \in j} \sum_{\mu \in \bar{\mathbb{F}}} d_{\mu\nu} \bar{x}_{\mu\nu} \\ &= \sum_{i \in \bar{\mathbb{F}}} d_{ij} \sum_{\nu \in j} \sum_{\mu \in i} \bar{x}_{\mu\nu} \\ &= \sum_{i \in \bar{\mathbb{F}}} d_{ij} x_{ij}^* = C_j^*.\end{aligned}$$

Thus, summing over all demands, the expected total connection cost is

$$\begin{aligned}\text{Exp}[C_{\text{EGUP}}] &\leq \sum_{j \in \mathbb{C}} \sum_{\nu \in j} (C_{\nu}^{\text{avg}} + 2\alpha_{\nu}^*) \\ &= \sum_{j \in \mathbb{C}} (C_j^* + 2r_j \alpha_j^*) = C^* + 2 \cdot \text{LP}^*,\end{aligned}$$

completing the proof of the lemma.  $\square$

**Theorem 8.** *Algorithm EGUP is a 3-approximation algorithm.*

*Proof.* By Property (SI.2), different demands from the same client are assigned to different primary demands, and by (PD.1) each primary demand opens a different facility. This ensures that our solution is feasible, namely each client  $j$  is connected to  $r_j$  different facilities (some possibly located on the same site). As for the total cost, Lemma 6 and Lemma 7 imply that the total cost is at most  $F^* + C^* + 2 \cdot \text{LP}^* = 3 \cdot \text{LP}^* \leq 3 \cdot \text{OPT}$ .  $\square$

## 6 Algorithm ECHS with Ratio 1.736

In this section we improve the approximation ratio to  $1 + 2/e \approx 1.736$ . The improvement comes from a slightly modified rounding process and refined analysis. Note that the facility opening cost of Algorithm EGUP does not exceed that of the fractional optimum solution, while the connection cost is quite far from the optimum, since we connect a non-primary demand to a facility in its assigned primary demand's neighborhood and then estimate the distance using the triangle inequality. To improve the estimate of the connection cost, the basic idea, following the approach of Chudak and Shmoys [5], is to connect a non-primary demand to its nearest neighbor when one is available and only use the facility that its primary demand opens when none of its neighbor is open.

**Algorithm ECHS.** As before, the algorithm starts by solving the linear program and applying the adaptive partitioning algorithm described in Section 4 to obtain a partitioned solution  $(\bar{x}, \bar{y})$ . Then we apply the rounding process to compute an integral solution (see Pseudocode 3). For convenience, we will use the term *facility cluster* for the neighborhood of a primary demand. Facilities that do not belong to these clusters are said to be *non-clustered*. We start, as before, by opening exactly one facility  $\phi(\kappa)$  in the facility cluster of each primary demand  $\kappa$  (Line 2). For any non-primary demand  $\nu$  assigned to  $\kappa$ , we refer to  $\phi(\kappa)$  as the *target* facility of  $\nu$ . In Algorithm EGUP,  $\nu$  was connected to  $\phi(\kappa)$ , but in Algorithm ECHS we may be able to find an open facility in  $\nu$ 's neighborhood and connect  $\nu$  to this facility. Specifically, the two changes in the algorithm are as follows:

- (1) Each non-clustered facility  $\mu$  is opened, independently, with probability  $\bar{y}_\mu$  (Lines 4–5). Notice that if  $\bar{y}_\mu > 0$  then, due to completeness of the partitioned fractional solution, we have  $\bar{y}_\mu = \bar{x}_{\mu\nu}$  for some demand  $\nu$ . This implies that  $\bar{y}_\mu \leq 1$ , because  $\bar{x}_{\mu\nu} \leq 1$ , by (PS.1).
- (2) When connecting demands to facilities, a primary demand  $\kappa$  is connected to the only facility  $\phi(\kappa)$  opened in its neighborhood, as before (Line 3). For a non-primary demand  $\nu$ , if its neighborhood  $\bar{N}(\nu)$  has an open facility, we connect  $\nu$  to the closest open facility in  $\bar{N}(\nu)$  (Line 8). Otherwise, we connect  $\nu$  to its target facility (Line 10).

---

### Pseudocode 3 Algorithm ECHS: Constructing Integral Solution

---

```

1: for each  $\kappa \in P$  do
2:   choose one  $\phi(\kappa) \in \bar{N}(\kappa)$ , with each  $\mu \in \bar{N}(\kappa)$  chosen as  $\phi(\kappa)$  with probability  $\bar{y}_\mu$ 
3:   open  $\phi(\kappa)$  and connect  $\kappa$  to  $\phi(\kappa)$ 
4: for each  $\mu \in \bar{\mathbb{F}} - \bigcup_{\kappa \in P} \bar{N}(\kappa)$  do
5:   open  $\mu$  with probability  $\bar{y}_\mu$  (independently)
6: for each non-primary demand  $\nu \in \bar{\mathbb{C}}$  do
7:   if any facility in  $\bar{N}(\nu)$  is open then
8:     connect  $\nu$  to the nearest open facility in  $\bar{N}(\nu)$ 
9:   else
10:    connect  $\nu$  to  $\phi(\kappa)$  where  $\kappa$  is  $\nu$ 's primary demand

```

---

**Analysis.** We shall first argue that the integral solution thus constructed is feasible, and then we bound the total cost of the solution. Regarding feasibility, the only constraint that is not explicitly

enforced by the algorithm is the fault-tolerance requirement; namely that each client  $j$  is connected to  $r_j$  different facilities. Let  $\nu$  and  $\nu'$  be two different sibling demands of client  $j$  and let their assigned primary demands be  $\kappa$  and  $\kappa'$  respectively. Due to (SI.2) we know  $\kappa \neq \kappa'$ . From (SI.1) we have  $\bar{N}(\nu) \cap \bar{N}(\nu') = \emptyset$ . From (SI.2), we have  $\bar{N}(\nu) \cap \bar{N}(\kappa') = \emptyset$  and  $\bar{N}(\nu') \cap \bar{N}(\kappa) = \emptyset$ . From (PD.1) we have  $\bar{N}(\kappa) \cap \bar{N}(\kappa') = \emptyset$ . It follows that  $(\bar{N}(\nu) \cup \bar{N}(\kappa)) \cap (\bar{N}(\nu') \cup \bar{N}(\kappa')) = \emptyset$ . Since the algorithm connects  $\nu$  to some facility in  $\bar{N}(\nu) \cup \bar{N}(\kappa)$  and  $\nu'$  to some facility in  $\bar{N}(\nu') \cup \bar{N}(\kappa')$ ,  $\nu$  and  $\nu'$  will be connected to different facilities.

We now show that the expected cost of the computed solution is bounded by  $(1 + 2/e) \cdot \text{LP}^*$ . By (PD.1), every facility may appear in at most one primary demand's neighborhood, and the facilities open in Line 4–5 of Pseudocode 3 do not appear in any primary demand's neighborhood. Therefore, by linearity of expectation, the expected facility cost of Algorithm ECHS is

$$\text{Exp}[F_{\text{ECHS}}] = \sum_{\mu \in \mathbb{F}} f_\mu \bar{y}_\mu = \sum_{i \in \mathbb{F}} f_i \sum_{\mu \in i} \bar{y}_\mu = \sum_{i \in \mathbb{F}} f_i y_i^* = F^*,$$

where the third equality follows from (PS.3).

To bound the connection cost, we adapt an argument of Chudak and Shmoys [5]. Consider a demand  $\nu$ . This demand can either get connected directly to some facility in  $\bar{N}(\nu)$  or indirectly to its target facility  $\phi(\kappa) \in \bar{N}(\kappa)$ , where  $\kappa$  is the primary demand to which  $\nu$  is assigned.

We now estimate the expected cost  $d_{\phi(\kappa)\nu}$  of the indirect connection. Let  $\Lambda_\nu$  denote the event that none of the facilities in  $\bar{N}(\nu)$  is opened. Then

$$\text{Exp}[d_{\phi(\kappa)\nu} | \Lambda_\nu] = D(\bar{N}(\kappa) \setminus \bar{N}(\nu), \nu),$$

where  $D(A, \sigma) \stackrel{\text{def}}{=} \sum_{\mu \in A} d_{\mu\sigma} \bar{y}_\mu / \sum_{\mu \in A} \bar{y}_\mu$ , for any set  $A$  of facilities and a demand  $\sigma$ . Note that  $C_\nu^{\text{avg}} = D(\bar{N}(\nu), \nu)$ , and that  $\Lambda_\nu$  implies that  $\bar{N}(\kappa) \setminus \bar{N}(\nu) \neq \emptyset$ , since  $\bar{N}(\kappa)$  contains at least one open facility, namely  $\phi(\kappa)$ .

**Lemma 9.** *Let  $\nu$  be a demand assigned to a primary demand  $\kappa$ , and assume that  $\bar{N}(\kappa) \setminus \bar{N}(\nu) \neq \emptyset$ . Then  $\text{Exp}[d_{\phi(\kappa)\nu} | \Lambda_\nu] \leq C_\nu^{\text{avg}} + 2\alpha_\nu^*$ .*

*Proof.* By the discussion above, we need to show that  $D(\bar{N}(\kappa) \setminus \bar{N}(\nu), \nu) \leq C_\nu^{\text{avg}} + 2\alpha_\nu^*$ . There are two cases to consider.

Case 1: There exists some  $\mu' \in \bar{N}(\kappa) \cap \bar{N}(\nu)$  such that  $d_{\mu'\kappa} \leq C_\kappa^{\text{avg}}$ . In this case, for every  $\mu \in \bar{N}(\kappa) \setminus \bar{N}(\nu)$ , we have

$$d_{\mu\nu} \leq d_{\mu\kappa} + d_{\mu'\kappa} + d_{\mu'\nu} \leq \alpha_\kappa^* + C_\kappa^{\text{avg}} + \alpha_\nu^* \leq C_\nu^{\text{avg}} + 2\alpha_\nu^*,$$

using the triangle inequality, complementary slackness, and (PD.3(b)). By summing over all  $\mu \in \bar{N}(\kappa) \setminus \bar{N}(\nu)$ , it follows that  $D(\bar{N}(\kappa) \setminus \bar{N}(\nu), \nu) \leq C_\nu^{\text{avg}} + 2\alpha_\nu^*$ .

Case 2: Every  $\mu' \in \bar{N}(\kappa) \cap \bar{N}(\nu)$  has  $d_{\mu'\kappa} > C_\kappa^{\text{avg}}$ . Then we have  $D(\bar{N}(\kappa) \setminus \bar{N}(\nu), \kappa) \leq D(\bar{N}(\kappa), \kappa) = C_\kappa^{\text{avg}}$ . Therefore, choosing an arbitrary  $\mu' \in \bar{N}(\kappa) \cap \bar{N}(\nu)$ , we obtain

$$D(\bar{N}(\kappa) \setminus \bar{N}(\nu), \nu) \leq D(\bar{N}(\kappa) \setminus \bar{N}(\nu), \kappa) + d_{\mu'\kappa} + d_{\mu'\nu} \leq C_\kappa^{\text{avg}} + \alpha_\kappa^* + \alpha_\nu^* \leq C_\nu^{\text{avg}} + 2\alpha_\nu^*,$$

where we again use the triangle inequality, complementary slackness, and (PD.3(b)).

Since the lemma holds in both cases, the proof is now complete.  $\square$

We now continue our estimation of the connection cost. We note first that for a primary demand  $\kappa$ , its expected connection cost is  $C_\kappa^{\text{avg}} = \sum_{\mu \in \bar{N}(\kappa)} d_{\mu\kappa} \bar{x}_{\mu\kappa}$  as in the previous section.

Next, we consider a non-primary demand  $\nu$ . Denote by  $C_\nu$  the random variable representing the connection cost for  $\nu$ . We first deal with the case when  $\bar{N}(\kappa) \setminus \bar{N}(\nu) \neq \emptyset$ , which is the same as  $\bar{N}(\kappa) \subseteq \bar{N}(\nu)$ . This actually implies that  $\bar{N}(\kappa) = \bar{N}(\nu)$  because Property (CO) implies that  $\bar{x}_{\mu\nu} = \bar{y}_\mu = \bar{x}_{\mu\kappa}$  for every  $\mu \in \bar{N}(\kappa)$  and therefore  $\sum_{\mu \in \bar{N}(\kappa)} \bar{x}_{\mu\nu} = 1$ . Thus  $\text{Exp}[C_\nu] = C_\nu^{\text{avg}}$ .

Now we bound the expected connection cost of  $\nu$  when  $\bar{N}(\kappa) \setminus \bar{N}(\nu) \neq \emptyset$ . Let  $\bar{N}(\nu) = \{\mu_1, \dots, \mu_l\}$  and let  $d_s = d_{\mu_s\nu}$  and  $y_s = \bar{y}_{\mu_s}$  for  $s = 1, \dots, l$ . By reordering, we can assume that  $d_1 \leq d_2 \leq \dots \leq d_l$ . By Pseudocode 3, the connection cost is no more than that obtained through the random process that opens each  $\mu_s$  independently with probability  $\bar{x}_{\mu_s\nu} = y_s$  (because  $\bar{x}_{\mu_s\nu} > 0$  and by (CO)), and connects  $\nu$  to the nearest such open facility, if any of them opens; otherwise  $\nu$  is connected indirectly to its target facility  $\phi(\kappa)$ . The intuition is that we only use a facility  $\mu_s$  if none of  $\mu_1, \dots, \mu_{s-1}$  is open. We know that  $\mu_s$  opens (unconditionally) with probability  $y_s$ . The only way that some of  $\mu_1, \dots, \mu_{s-1}$  can affect that probability is if they belong to  $\bar{N}(\kappa)$  for some primary demand  $\kappa$  and it also happens to be that  $\mu_s \in \bar{N}(\kappa)$  as well. However in this case, the condition that they are closed actually implies that the conditional probability of  $\mu_s$  being open is larger than  $y_s$ . (For a detailed proof, see [5].)

Putting it all together, we estimate the (unconditional) expected connection cost of  $\nu$  as follows:

$$\begin{aligned} \text{Exp}[C_\nu] &\leq \sum_{r=1}^l d_r y_r \prod_{s=1}^{r-1} (1 - y_s) + \text{Exp}[d_{\phi(\kappa)\nu} | \Lambda_\nu] \cdot \prod_{s=1}^l (1 - y_s) \\ &\leq \left(1 - \prod_{s=1}^l (1 - y_s)\right) \cdot \sum_{r=1}^l d_r y_r + \text{Exp}[d_{\phi(\kappa)\nu} | \Lambda_\nu] \cdot \prod_{s=1}^l (1 - y_s) \\ &\leq \left(1 - \frac{1}{e}\right) \sum_{r=1}^l d_r y_r + \frac{1}{e} \text{Exp}[d_{\phi(\kappa)\nu} | \Lambda_\nu] \leq \left(1 - \frac{1}{e}\right) C_\nu^{\text{avg}} + \frac{1}{e} (C_\nu^{\text{avg}} + 2\alpha_\nu^*) = C_\nu^{\text{avg}} + \frac{2}{e} \alpha_\nu^*, \end{aligned}$$

where the second inequality is shown in the appendix (see also [5]) and the last inequality follows from Lemma 9.

Summing over all demands of a client  $j$ , we bound the expected connection cost of client  $j$ :

$$\text{Exp}[C_j] \leq \sum_{\nu \in j} (C_\nu^{\text{avg}} + \frac{2}{e} \alpha_\nu^*) = C_j^* + \frac{2}{e} r_j \alpha_j^*.$$

Finally, summing over all clients  $j$ , we obtain our bound on the expected connection cost,

$$\text{Exp}[C_{\text{ECHS}}] \leq C^* + \frac{2}{e} \text{LP}^*.$$

Therefore, we have established that our algorithm constructs a feasible integral solution with an overall expected cost

$$\text{Exp}[F_{\text{ECHS}} + C_{\text{ECHS}}] \leq F^* + C^* + \frac{2}{e} \cdot \text{LP}^* = (1 + 2/e) \cdot \text{LP}^* \leq (1 + 2/e) \cdot \text{OPT}.$$

Summarizing, we obtain the main result of this section.

**Theorem 10.** *Algorithm ECHS is a  $(1 + 2/e)$ -approximation algorithm for FTFP.*



## 7 Algorithm EBGs with Ratio 1.575

In this section we give our main result, a 1.575-approximation algorithm for FTFP, where 1.575 is the value of  $\min_{\gamma \geq 1} \max\{\gamma, 1 + 2/e^\gamma, \frac{1/e + 1/e^\gamma}{1 - 1/\gamma}\}$ , rounded to three decimal digits. This matches the ratio of the best known LP-rounding algorithm for UFL by Byrka *et al.* [3]. Recall that in Section 6 we showed how to compute an integral solution with facility cost bounded by  $F^*$  and connection cost bounded by  $C^* + 2/e \cdot \text{LP}^*$ . A natural idea is to balance these two costs, by reducing the connection cost, at the expense of slightly increasing the facility cost. If it works, this should result in reducing the approximation ratio.

Our approach can be thought of as a combination of the ideas in [3] with the techniques of demand reduction and adaptive partitioning that we introduced earlier. However, our adaptive partitioning technique needs to be carefully modified because now we will be using a more intricate neighborhood structure, with the neighborhood of each demand divided into two parts, the close and far neighborhood, and with some conditions on which pairs of neighborhoods need to overlap and which need to be disjoint. The final rounding stage that construct an integral solution is a relatively straightforward generalization of the rounding method in [3].

We begin by describing properties that our partitioned fractional solution  $(\bar{x}, \bar{y})$  needs to satisfy. The neighborhood  $\bar{N}(\nu)$  of each demand  $\nu$  will be divided into two disjoint parts. The first part, called the *close neighborhood*  $\bar{N}_{\text{cls}}(\nu)$ , contains the facilities in  $\bar{N}(\nu)$  nearest to  $\nu$  with the total connection value equal  $1/\gamma$ . The second part, called the *far neighborhood*  $\bar{N}_{\text{far}}(\nu)$ , contains the remaining facilities in  $\bar{N}(\nu)$ . The formal definitions of these sets are given below in Property (NB). The respective average connection costs from  $\nu$  for these sets are defined by  $C_{\text{cls}}^{\text{avg}}(\nu) = \gamma \sum_{\mu \in \bar{N}_{\text{cls}}(\nu)} d_{\mu\nu} \bar{x}_{\mu\nu}$  and  $C_{\text{far}}^{\text{avg}}(\nu) = \frac{\gamma}{\gamma-1} \sum_{\mu \in \bar{N}_{\text{far}}(\nu)} d_{\mu\nu} \bar{x}_{\mu\nu}$ . We will also use notation  $C_{\text{cls}}^{\text{max}}(\nu) = \max_{\mu \in \bar{N}_{\text{cls}}(\nu)} d_{\mu\nu}$  for the maximum distance from  $\nu$  to its close neighborhood.

Our partitioned solution  $(\bar{x}, \bar{y})$  must satisfy the same partitioning and completeness properties as before, namely properties (PS) and (CO) in Section 4. In addition, it must satisfy new neighborhood property (NB) and modified properties (PD) and (SI), listed below.

(NB) For each demand  $\nu$ , its neighborhood is divided into *close* and *far* neighborhood, that is  $\bar{N}(\nu) = \bar{N}_{\text{cls}}(\nu) \cup \bar{N}_{\text{far}}(\nu)$ , where

- $\bar{N}_{\text{cls}}(\nu) \cap \bar{N}_{\text{far}}(\nu) = \emptyset$ ,
- $\sum_{\mu \in \bar{N}_{\text{cls}}(\nu)} \bar{x}_{\mu\nu} = 1/\gamma$ , and
- if  $\mu \in \bar{N}_{\text{cls}}(\nu)$  and  $\mu' \in \bar{N}_{\text{far}}(\nu)$  then  $d_{\mu\nu} \leq d_{\mu'\nu}$ .

Note that the second condition, together with (PS.1), implies that  $\sum_{\mu \in \bar{N}_{\text{far}}(\nu)} \bar{x}_{\mu\nu} = 1 - 1/\gamma$ .

(PD') *Primary demands*. Primary demands satisfy the following conditions:

1. For any two different primary demands  $\kappa, \kappa' \in P$  we have  $\bar{N}_{\text{cls}}(\kappa) \cap \bar{N}_{\text{cls}}(\kappa') = \emptyset$ .
2. For each site  $i \in \mathbb{F}$ ,  $\sum_{\kappa \in P} \sum_{\mu \in i \cap \bar{N}_{\text{cls}}(\kappa)} \bar{x}_{\mu\kappa} \leq y_i^*$ .
3. Each demand  $\nu \in \bar{\mathbb{C}}$  is assigned to one primary demand  $\kappa \in P$  such that
  - (a)  $\bar{N}_{\text{cls}}(\nu) \cap \bar{N}_{\text{cls}}(\kappa) \neq \emptyset$ , and
  - (b)  $C_{\text{cls}}^{\text{avg}}(\nu) + C_{\text{cls}}^{\text{max}}(\nu) \geq C_{\text{cls}}^{\text{avg}}(\kappa) + C_{\text{cls}}^{\text{max}}(\kappa)$ .

(SI') *Siblings*. For any pair  $\nu, \nu'$  of different siblings we have

1.  $\bar{N}(\nu) \cap \bar{N}(\nu') = \emptyset$ .
2. If  $\nu$  is assigned to a primary demand  $\kappa$  then  $\bar{N}(\nu') \cap \bar{N}_{\text{cls}}(\kappa) = \emptyset$ . In particular, by Property (PD.3(a)), this implies that different sibling demands are assigned to different primary demands.

**Modified adaptive partitioning.** To obtain a fractional solution with the above properties, we employ a modified adaptive partitioning algorithm. As in Section 4, we have two phases. In Phase 1 we split clients into demands and create facilities on sites, while in Phase 2 we augment each demand's connection values so that its total value is 1.

Phase 1 runs in iterations. Consider any client  $j$ . As before,  $\tilde{N}(j)$  is the neighborhood of  $j$  with respect to the yet unpartitioned solution, namely the set of facilities  $\mu$  such that  $\tilde{x}_{\mu j} > 0$ . Order the facilities in this set as  $\tilde{N}(j) = \{\mu_1, \dots, \mu_q\}$  in order of non-decreasing distance from  $j$ , that is  $d_{\mu_1 j} \leq d_{\mu_2 j} \leq \dots \leq d_{\mu_q j}$ , where  $q = |\tilde{N}(j)|$ . Without loss of generality, there is an index  $l$  for which  $\sum_{s=1}^l \tilde{x}_{s j} = 1/\gamma$ , since we can always split one facility to have this property. Then we define  $\tilde{N}_\gamma(j) = \{\mu_1, \dots, \mu_l\}$ . We also use notation

$$\text{tcc}_\gamma(j) = D(\tilde{N}_\gamma(j), j) = \sum_{\mu \in \tilde{N}_\gamma(j)} d_{\mu j} \tilde{x}_{\mu j} \quad \text{and} \quad \text{dmax}_\gamma(j) = \max_{\mu \in \tilde{N}_\gamma(j)} d_{\mu j}.$$

In each iteration, we find a not yet exhausted client  $p$  that minimizes the value of  $\text{tcc}_\gamma(p) + \text{dmax}_\gamma(p)$ . Now we have two cases:

Case 1:  $\tilde{N}_\gamma(p) \cap \bar{N}_{\text{cls}}(\kappa) \neq \emptyset$ , for some existing primary demand  $\kappa$ . In this case we assign  $\nu$  to  $\kappa$ . As before, if there are multiple such  $\kappa$ , we pick any of them. We also fix  $\bar{x}_{\mu\nu} \leftarrow \tilde{x}_{\mu p}$ ,  $\tilde{x}_{\mu p} \leftarrow 0$  for each  $\mu \in \tilde{N}_\gamma(p) \cap \bar{N}_{\text{cls}}(\kappa)$ . As before, although we check for overlap between  $\tilde{N}_\gamma(p)$  and  $\bar{N}_{\text{cls}}(\kappa)$ , the facilities we actually move into  $\bar{N}(\nu)$  include all facilities in the intersection of  $\tilde{N}(p)$ , a bigger set, with  $\bar{N}_{\text{cls}}(\kappa)$ . We would like to point out that  $\bar{N}(\nu)$  is not finalized at this time as we will add more facilities to it in the augment phase. As a result  $\bar{N}_{\text{cls}}(\nu)$  is not fixed either, as we could potentially add facilities closer to  $\nu$  than facilities already in  $\bar{N}(\nu)$ . Recall that by definition  $\bar{N}_{\text{cls}}(\nu)$  consists of the facilities that are closest to  $\nu$  once  $\bar{N}(\nu)$  is fixed with total connection value of 1.

Case 2:  $\tilde{N}_\gamma(p) \cap \bar{N}_{\text{cls}}(\kappa) = \emptyset$ , for all existing primary demands  $\kappa$ . In this case we make  $\nu$  a primary demand. We then fix  $\bar{x}_{\mu\nu} \leftarrow \tilde{x}_{\mu p}$  for  $\mu \in \tilde{N}_\gamma(p)$  and set the corresponding  $\tilde{x}_{\mu p}$  to 0. Note that the total connection value in  $\bar{N}_{\text{cls}}(\nu)$  is now exactly  $1/\gamma$ . The set  $\tilde{N}_\gamma(p)$  turns out to coincide with  $\bar{N}_{\text{cls}}(\nu)$  as the facilities in  $\tilde{N}(p) \setminus \tilde{N}_\gamma(p)$  are all farther away than any facility in  $\tilde{N}_\gamma(p)$ . In the augment phase we have available only facilities in some subset of  $\tilde{N}(p) \setminus \tilde{N}_\gamma(p)$ . Thus  $\bar{N}_{\text{cls}}(\nu)$  is defined when it is created.

Once all clients are exhausted, that is, each client  $j$  has  $r_j$  demands created, Phase 1 concludes. We then do Phase 2, the augmenting phase. For each demand  $\nu$  of client  $j$  with total connection value less than 1, we use our `AUGMENTTOUNIT()` procedure to add additional facilities from  $\tilde{N}(j)$  to its neighborhood to make its total connection value equal 1, as before a facility is removed from

$\tilde{N}(j)$  once added to a demand's neighborhood. We do facility split if necessary to make  $\bar{N}(\nu)$  have total connection value of 1. This completes the description of the partitioning algorithm.

We now argue that the fractional solution  $(\bar{x}, \bar{y})$  satisfies all the stated properties. Properties (PS), (CO), (PD.1) and (SI.1) are directly enforced by the adaptive partitioning algorithm. The proofs for other properties are similar to those in Section 4, with the exception of (PD.3(a)), which we justify below.

The argument for (PD.3(a)) is a bit subtle, because of possible complications arising in the augmentation phase. This phase does not change close neighborhoods of primary demands, as each primary demand already contains all the nearest facilities with total connection value  $1/\gamma$ . For non-primary demands, however,  $\bar{N}(\nu)$ , for  $\nu \in j$ , takes all facilities in  $\bar{N}_{\text{cls}}(\kappa) \cap \tilde{N}(j)$ , which might be close to  $\kappa$  but far from  $j$ . It seems that facilities added in the augment step might actually be closer to  $\nu$  than some of the facilities already in  $\bar{N}(\nu)$ . As a result, facilities added in the augment step might appear in  $\bar{N}_{\text{cls}}(\nu)$ , yet they are not in  $\bar{N}_{\text{cls}}(\kappa)$ , the close neighborhood of the primary demand  $\kappa$  that  $\nu$  is assigned to. Nevertheless, we show that Property (PD.3(a)) holds.

Consider an iteration when we create a demand  $\nu \in p$  and assign it to  $\kappa$ . Then the set  $B(p) = \tilde{N}_\gamma(p) \cap \bar{N}_{\text{cls}}(\kappa)$  is not empty. We claim that  $B(p)$  must be a subset of  $\bar{N}_{\text{cls}}(\nu)$  after  $\bar{N}(\nu)$  is finalized with a total connection value of 1. To see this, first observe that  $B(p)$  is a subset of  $\bar{N}(\nu)$ , which in turn is a subset of  $\tilde{N}(p)$ , after taking into account the facility split. Here  $\tilde{N}(p)$  refers to the neighborhood of client  $p$  just before  $\nu$  was created. For an arbitrary set of facilities  $A$  define  $\text{dmax}(A, \nu)$  as the minimum distance  $\tau$  such that  $\sum_{\mu \in A: d_{\mu\nu} \leq \tau} \bar{y}_\mu \geq 1/\gamma$ . Adding additional facilities into  $A$  cannot make  $\text{dmax}(A, \nu)$  larger, so it follows that  $\text{dmax}(\bar{N}_{\text{cls}}(\nu), \nu) \geq \text{dmax}(\tilde{N}(p), \nu)$ , because  $\bar{N}_{\text{cls}}(\nu)$  is a subset of  $\tilde{N}(p)$ . Since we have  $d_{\mu\nu} = d_{\mu p}$  by definition, it is easy to see that every  $\mu \in B(p)$  satisfies  $d_{\mu\nu} \leq \text{dmax}(\tilde{N}(p), \nu) \leq \text{dmax}(\bar{N}_{\text{cls}}(\nu), \nu)$  and hence they all belong to  $\bar{N}_{\text{cls}}(\nu)$ . We need to be a bit more careful here when we have a tie in  $d_{\mu\nu}$  but we can assume ties are always broken in favor of facilities in  $B(p)$  when defining  $\bar{N}_{\text{cls}}(\nu)$ . Finally, since  $B(p) \neq \emptyset$ , we have that the close neighborhood of a demand  $\nu$  and its primary demand  $\kappa$  must overlap.

**Algorithm EBGs.** The complete algorithm starts with solving the linear program and computing the partitioning described earlier in this section. Given the partitioned fractional solution  $(\bar{x}, \bar{y})$  with the desired properties, we then start opening facilities and making connections to obtain an integral solution. As before, we open exactly one facility in each cluster (the close neighborhood of a primary demand), but now each facility  $\mu$  is chosen with probability  $\gamma \bar{y}_\mu$ . The non-clustered facilities  $\mu$ , those that do not belong to  $\bar{N}_{\text{cls}}(\kappa)$  for any primary demand  $\kappa$ , are opened independently with probability  $\gamma \bar{y}_\mu$  each.

Next, we connect demands to facilities. Each primary demand  $\kappa$  will connect to the only facility  $\phi(\kappa)$  open in its cluster  $\bar{N}_{\text{cls}}(\kappa)$ . For each non-primary demand  $\nu$ , if there is an open facility in  $\bar{N}(\nu)$  then we connect  $\nu$  to the nearest such facility. Otherwise, we connect  $\nu$  to its *target facility*  $\phi(\kappa)$ , where  $\kappa$  is the primary demand that  $\nu$  is assigned to.

**Analysis.** The feasibility of our integral solution follows from Properties (SI.1), (SI.2), and (PD.1), as these properties together ensure that each facility is accessible to at most one demand among sibling demands of the same client, regardless whether a demand connects to its neighbor or its target facility.

The expected facility cost of our algorithm is bounded by  $\gamma F^*$ , using essentially the same argument as in the previous section (with the factor  $\gamma$  accounting for using probabilities  $\gamma \bar{y}_\mu$  instead of  $\bar{y}_\mu$ ).

We now bound the connection cost. Properties (PD.3(a)) and (PD.3(b)) allow us to bound the expected distance from a demand  $\nu$  to its target facility by  $C_{\text{cls}}^{\text{avg}}(\nu) + C_{\text{cls}}^{\text{max}}(\nu) + C_{\text{far}}^{\text{avg}}(\nu)$ , in the event that none of  $\nu$ 's neighbors opens, using a similar argument as Lemma 2.2 in [3]<sup>2</sup>. We are then able to estimate the expected connection cost for demand  $\nu$  using an argument similar to [3]: with probability no less than  $1 - 1/e$ ,  $\nu$  has some facility open in its close neighborhood, with probability no less than  $1 - 1/e^\gamma$ ,  $\nu$  has some facility open in its overall neighborhood, and with probability no more than  $1/e^\gamma$ ,  $\nu$  will connect to its target facility. This gives us the bound

$$\begin{aligned} \text{Exp}[C_\nu] &\leq C_{\text{cls}}^{\text{avg}}(\nu)(1 - 1/e) + C_{\text{far}}^{\text{avg}}(\nu)(1/e - 1/e^\gamma) + (C_{\text{cls}}^{\text{avg}}(\nu) + C_{\text{cls}}^{\text{max}}(\nu) + C_{\text{far}}^{\text{avg}}(\nu))1/e^\gamma \\ &\leq C_{\text{cls}}^{\text{avg}}(\nu)(1 - 1/e) + C_{\text{far}}^{\text{avg}}(\nu)(1/e - 1/e^\gamma) + (C_{\text{cls}}^{\text{avg}}(\nu) + 2C_{\text{far}}^{\text{avg}}(\nu))1/e^\gamma \\ &\leq C^{\text{avg}}(\nu)((1 - \rho_\nu)(\frac{1/e + 1/e^\gamma}{1 - 1/\gamma}) + \rho_\nu(1 + 2/e^\gamma)) \\ &\leq C^{\text{avg}}(\nu) \cdot \max\left\{\frac{1/e + 1/e^\gamma}{1 - 1/\gamma}, 1 + \frac{2}{e^\gamma}\right\}, \end{aligned}$$

where  $\rho_\nu = C_{\text{cls}}^{\text{avg}}(\nu)/C^{\text{avg}}(\nu)$ . It is easy to see that  $\rho_\nu$  is between 0 and 1. Since  $\sum_{\nu \in j} C^{\text{avg}}(\nu) = \sum_{\nu \in j} \sum_{\mu \in \mathbb{F}} d_{\mu\nu} \bar{x}_{\mu\nu} = \sum_{i \in \mathbb{F}} d_{ij} x_{ij}^* = C_j^*$ , summing over all clients  $j$  we have total connection cost bounded by  $C^* \max\{\frac{1/e + 1/e^\gamma}{1 - 1/\gamma}, 1 + \frac{2}{e^\gamma}\}$ . The expected facility cost is bounded by  $\gamma F^*$ , as argued earlier. Hence the total cost is bounded by  $\max\{\gamma, \frac{1/e + 1/e^\gamma}{1 - 1/\gamma}, 1 + \frac{2}{e^\gamma}\} \cdot \text{LP}^*$ . Picking  $\gamma = 1.575$  we obtain the desired ratio.

**Theorem 11.** *Algorithm EBGs is a 1.575-approximation algorithm for FTFP.*

## 8 Final Comments

In this paper we show a sequence of LP-rounding approximation algorithms for FTFP, with the best algorithm achieving ratio 1.575. The two techniques we introduced, namely the demand reduction and adaptive partitioning, are very flexible. Should any new LP-rounding algorithms be discovered for UFL, we believe that with our approach they can be adapted to FTFP as well, preserving the approximation ratio. In fact, by randomizing the scaling parameter  $\gamma$  from Section 7, following the approach by Li [13], we could further improve the ratio to below 1.575. This is not enough, however, to match the 1.488 bound for UFL in [13], because matching this bound also requires appropriately extending dual-fitting algorithms [15] to FTFP, which we have so far been unable to do.

One of the main open problems in this area is whether FTFL can be approximated with the same ratio as UFL, and our work was partly motivated by this question. The techniques we introduced are not directly applicable to FTFL, mainly because our partitioning approach involves facility splitting that could result in several sibling demands being served by facilities on the same site. Nonetheless, we hope that further refinements of our construction might get around this issue and lead to new algorithms for FTFL with improved ratios.

<sup>2</sup>The full proof of the lemma appears in [2] as Lemma 3.3.

## References

- [1] Vijay Arya, Naveen Garg, Rohit Khandekar, Adam Meyerson, Kamesh Munagala, and Vinayaka Pandit. Local search heuristics for k-median and facility location problems. *SIAM J. Comput.*, 33(3):544–562, 2004.
- [2] Jaroslaw Byrka and Karen Aardal. An optimal bifactor approximation algorithm for the metric uncapacitated facility location problem. *SIAM J. Comput.*, 39(6):2212–2231, 2010.
- [3] Jaroslaw Byrka, MohammadReza Ghodsi, and Aravind Srinivasan. LP-rounding algorithms for facility-location problems. *CoRR*, abs/1007.3611, 2010.
- [4] Jaroslaw Byrka, Aravind Srinivasan, and Chaitanya Swamy. Fault-tolerant facility location, a randomized dependent LP-rounding algorithm. In *Proceedings of the 14th Integer Programming and Combinatorial Optimization, IPCO '10*, pages 244–257, 2010.
- [5] Fabián Chudak and David Shmoys. Improved approximation algorithms for the uncapacitated facility location problem. *SIAM J. Comput.*, 33(1):1–25, 2004.
- [6] Sudipto Guha and Samir Khuller. Greedy strikes back: improved facility location algorithms. In *Proceedings of the 9th ACM-SIAM Symposium on Discrete Algorithms, SODA '98*, pages 649–657, 1998.
- [7] Sudipto Guha, Adam Meyerson, and Kamesh Munagala. Improved algorithms for fault tolerant facility location. In *Proceedings of the 12th ACM-SIAM Symposium on Discrete Algorithms, SODA '01*, pages 636–641, 2001.
- [8] Anupam Gupta. Lecture notes: CMU 15-854b, spring 2008, 2008.
- [9] Godfrey Harold Hardy, John Edensor Littlewood, and George Pólya. *Inequalities*. Cambridge University Press, 1988.
- [10] Kamal Jain, Mohammad Mahdian, Evangelos Markakis, Amin Saberi, and Vijay Vazirani. Greedy facility location algorithms analyzed using dual fitting with factor-revealing LP. *J. ACM*, 50(6):795–824, 2003.
- [11] Kamal Jain and Vijay Vazirani. Approximation algorithms for metric facility location and k-median problems using the primal-dual schema and lagrangian relaxation. *J. ACM*, 48(2):274–296, 2001.
- [12] Kamal Jain and Vijay V. Vazirani. An approximation algorithm for the fault tolerant metric facility location problem. *Algorithmica*, 38(3):433–439, 2003.
- [13] Shi Li. A 1.488 approximation algorithm for the uncapacitated facility location problem. In *Proceedings of the 38th International Conference on Automata, Languages and Programming, ICALP '11*, volume 6756, pages 77–88, 2011.
- [14] Kewen Liao and Hong Shen. Unconstrained and constrained fault-tolerant resource allocation. In *Proceedings of the 17th Annual International Conference on Computing and Combinatorics, COCOON'11*, pages 555–566, 2011.

- [15] Mohammad Mahdian, Yinyu Ye, and Jiawei Zhang. Approximation algorithms for metric facility location problems. *SIAM J. Comput.*, 36(2):411–432, 2006.
- [16] David Shmoys, Éva Tardos, and Karen Aardal. Approximation algorithms for facility location problems (extended abstract). In *Proceedings of the 29th Annual ACM Symposium on Theory of Computing, STOC '97*, pages 265–274, 1997.
- [17] Maxim Sviridenko. An improved approximation algorithm for the metric uncapacitated facility location problem. In *Proceedings of the 9th International IPCO Conference on Integer Programming and Combinatorial Optimization, IPCO '02*, pages 240–257, 2002.
- [18] Chaitanya Swamy and David Shmoys. Fault-tolerant facility location. *ACM Trans. Algorithms*, 4(4):1–27, 2008.
- [19] Jens Vygen. *Approximation Algorithms for Facility Location Problems*. Forschungsinst. für Diskrete Mathematik, 2005.
- [20] Shihong Xu and Hong Shen. The fault-tolerant facility allocation problem. In *Proceedings of the 20th International Symposium on Algorithms and Computation, ISAAC '09*, pages 689–698, 2009.
- [21] Li Yan and Marek Chrobak. Approximation algorithms for the fault-tolerant facility placement problem. *Inf. Process. Lett.*, 111(11):545–549, 2011.

## A An Elementary Proof of the Expected Connection Cost

In the  $1 + 2/e = 1.736$ -approximation in Section 6 we need to show the following inequality

$$\sum_{r=1}^l d_r y_r \prod_{s=1}^{r-1} (1 - y_s) \leq \left(1 - \prod_{s=1}^l (1 - y_s)\right) \cdot \sum_{r=1}^l d_r y_r \quad (3)$$

for  $d_1 \leq d_2 \leq \dots \leq d_l$  and  $\sum_{s=1}^l y_s = 1, y_s \geq 0$ .

In this section we give a new proof of this inequality, much simpler than the existing proof in [5], and also simpler than the argument by Sviridenko [17]. We derive this inequality from the following generalized version of the Chebyshev Sum Inequality:

$$\sum_i p_i \sum_j p_j a_j b_j \leq \sum_i p_i a_i \sum_j p_j b_j, \quad (4)$$

where each summation below runs from 1 to  $l$  and the sequences  $(a_i)$ ,  $(b_i)$  and  $(p_i)$  satisfy the following conditions:  $p_i \geq 0, a_i \geq 0, b_i \geq 0$  for all  $i$ ,  $a_1 \leq a_2 \leq \dots \leq a_l$ , and  $b_1 \geq b_2 \geq \dots \geq b_l$ .

Given inequality (4), we can obtain our inequality (3) by simple substitution

$$p_i \leftarrow y_i, a_i \leftarrow d_i, b_i \leftarrow \prod_{s=1}^{i-1} (1 - y_s)$$

For the sake of completeness, we include the proof of inequality (4), due to Hardy, Littlewood

and Polya [9]. The idea is to evaluate the following sum:

$$\begin{aligned}
S &= \sum_i p_i \sum_j p_j a_j b_j - \sum_i p_i a_i \sum_j p_j b_j \\
&= \sum_i \sum_j p_i p_j a_j b_j - \sum_i \sum_j p_i a_i p_j b_j \\
&= \sum_j \sum_i p_j p_i a_i b_i - \sum_j \sum_i p_j a_j p_i b_i \\
&= \frac{1}{2} \cdot \sum_i \sum_j (p_i p_j a_j b_j - p_i a_i p_j b_j + p_j p_i a_i b_i - p_j a_j p_i b_i) \\
&= \frac{1}{2} \cdot \sum_i \sum_j p_i p_j (a_i - a_j)(b_i - b_j) \leq 0.
\end{aligned}$$

The last inequality holds because  $(a_i - a_j)(b_i - b_j) \leq 0$ , since the sequences  $(a_i)$  and  $(b_i)$  are ordered oppositely.