

UNIVERSITY OF CALIFORNIA
RIVERSIDE

Approximation Algorithms for the Fault-Tolerant Facility Placement Problem

A Dissertation submitted in partial satisfaction
of the requirements for the degree of

Doctor of Philosophy

in

Computer Science

by

Li Yan

June 2013

Dissertation Committee:

Professor Marek Chrobak, Chairperson
Professor Tao Jiang
Professor Stefano Lonardi
Professor Neal Young

Copyright by
Li Yan
2013

The Dissertation of Li Yan is approved:

Committee Chairperson

University of California, Riverside

Acknowledgments

I would thank my advisor, Marek Chrobak, for bringing me into the PhD program of University of California Riverside, and for his guidance and patience on my study and research in the past five years. I am also grateful for the committee, Tao Jiang, Stefano Lonardi, and Neal Young for helpful discussions and comments on my research and this dissertation.

The supportive environment of the algorithm lab and computer science department has made PhD study here a pleasant experience and I am grateful for Claire Huang, Wei Li and the algorithm lab for helpful discussions and stimulation of ideas.

To my parents, who always have faith on my endeavor.

ABSTRACT OF THE DISSERTATION

Approximation Algorithms for the Fault-Tolerant Facility Placement Problem

by

Li Yan

Doctor of Philosophy, Graduate Program in Computer Science
University of California, Riverside, June 2013
Professor Marek Chrobak, Chairperson

In this thesis we study the Fault-Tolerant Facility Placement problem (FTFP). In the FTFP problem, we are given a set of sites where we can open facilities, and a set of clients with demands. To satisfy their demands, clients need to be connected to open facilities in sites. The goal is to satisfy all clients' demands while minimizing the total cost, that is the cost of opening facilities and the cost of connecting clients to facilities. The problem is shown to be NP-hard and hence we study approximation algorithms and their performance guarantee. Approximation algorithms are algorithms that run in polynomial time with provable performance when compared to optimal solutions.

We present two techniques that lead to several LP-rounding algorithms with progressively improved approximation ratio. The best ratio we have is 1.575. We also study the primal-dual approaches. In particular, we show that a natural greedy algorithm analyzed using the dual-fitting technique gives an approximation ratio of $O(\log n)$. On the negative side, under a natural assumption, we give an example showing that the dual-fitting analysis cannot give a ratio better than $O(\log n / \log \log n)$.

Contents

List of Figures	x
List of Tables	xi
1 Introduction	1
1.1 The Problem and the Background	1
1.2 Hardness Results on Approximating UFL and FTFP	6
2 Related Work and Known Results	12
2.1 Related Work on UFL	14
2.2 Related Work on FTFL	27
2.3 Our Problem FTFP	28
3 Linear Program	29
3.1 Notation and Definition	29
4 Techniques	34
4.1 Demand Reduction	35
4.1.1 Reduction from FTFP to FTFL	37
4.1.2 Asymptotic Approximation Ratio for Large Demands	38
4.2 Adaptive Partition	39
5 LP-rounding Algorithms	54
5.1 Algorithm EGUP with Ratio 3	54
5.2 Algorithm ECHS with Ratio 1.736	57
5.3 Algorithm EBGs with Ratio 1.575	65
6 Primal-dual Algorithms	86
6.1 The Greedy algorithm with $O(\log n)$ Ratio	90
6.2 An Example Showing the Difficulty in Obtaining $O(1)$ Ratio	94
7 Conclusion	100

Bibliography	102
A Technical Background	105
A.1 Linear Programming and Integer Programming	105
A.2 Proof of Inequality (5.3)	109

List of Figures

5.1	Illustration of the sets in the proof of Lemma 21	73
6.1	$\Omega(\log n / \log \log n)$ example for dual-fitting on FTFP	96

List of Tables

1.1	A history of approximation algorithms and their ratio for the UFL problem.	5
4.1	An example of an execution of the partitioning algorithm.	46

Chapter 1

Introduction

1.1 The Problem and the Background

In a broad sense, facility location problems are about selecting, from a set of sites, to open facilities, and servicing clients by connecting them to facilities. A classical application is to set up warehouses to distribute commodities to retailers. We want to have more warehouses so that every retailer can be close to some warehouse and thus save on shipping cost. On the other hand, setting up and maintaining a warehouse can be costly. So we would have to face the trade-off between having more warehouses and cheaper shipping, and having fewer warehouses for the sake of warehouse setup and maintenance cost, at the expense of shipping. Another application is to place content servers in a computer network. In the network we have a number of client machines that need to access files from one of the servers. Having more servers allows faster access for all clients. However, keeping a large number of servers around requires both substantial hardware and software investment, as

well as committing operation engineers to keep the servers up and running. A more recent example concerns today's web giants like google, facebook and amazon. Those web-based companies all have geographically distributed development offices. All these offices require access of large amount of data from a handful of data centers, located in a few carefully chosen locations. It is out of question that building a data center is costly. The electricity bill alone would discourage the plan of having an excessive number of data centers around. On the other hand, demands from all offices need to be addressed or engineers would be idling while waiting for data transfer to complete. We shall keep using the development office and data center example to illustrate several aspects of facility location problems.

First, setting up a data center incurs cost, and the cost varies between different locations. The real estate rent in Kansas usually do not compare to even a fraction of that in San Francisco or New York. The electricity rate is also very different in midwest area and California. Other factors like potential natural disasters like tonardo or earthquake could also be factored into the cost estimate to build a data center in a certain location. In short, different location has a different cost to set up a facility.

Second, different offices may have different demands. A larger office may require access to several data centers, either for speedy data transfer, or for redundancy when one of the data centers goes down.

Third, data centers may have capacities. It is thus reasonable to put a cap on the number of offices that a certain data center is able to serve. However, to keep the problem simple, we may or may not have this constraint in our problem definition.

Facility location problems have long been an active topic in both Operations Re-

search and Computer Science research communities. Problems that have been considered include

- the Uncapacitated Facility Location problem, where each site can open at most one facility and each client needs to be connected to one open facility,
- the Fault-Tolerant Facility Location problem, where each site can open at most one facility, and each client has a demand, which is the number of facilities that it needs to connect to,
- the Capacitated Facility Location problem, where a site can open at most one facility and each client connects to one facility, but every facility has a capacity, which is the maximum number of clients it can accept,
- the Prize-Collecting Facility Location problem, where a site can open at most one facility and each client can either connect to one facility, or stay unconnected, in which case a penalty is counted towards final cost.

In all the problems above, there is a cost to open a facility at a site, and there is a cost to connect a client to a facility as well. The goal is to minimize the total cost of opening facilities and connecting clients to facilities. In the Prize-Collecting Facility Location problem we also pay a penalty for each client not connected.

The UFL Problem The simplest variant, known as the uncapacitated facility location problem (UFL), is also the one that has been studied most extensively. A history of results and techniques for the UFL problem is given in Table 1.1. In the UFL problem, we are given a set of sites and a set of clients. Each site can open one facility, or none. Each client

has a demand of 1, meaning it needs to be connected to an open facility. Facilities are uncapacitated, meaning a facility can accept any number of connections from clients. We are also given the cost of opening one facility at each site, and the distance between a site and a client. The problem asks us to find a set of sites to open facilities, and to connect every client to an open facility, in a way that the total cost is minimized. The total cost is the sum of the cost to open facilities and the cost to connect clients to facilities.

The UFL problem with general distances has an algorithm with approximation ratio $O(\log n)$ where n is the number of clients, by Hochbaum [17]. A matching lower bound of $\Omega(\log n)$ is immediate, as the UFL problem contains the well-known Set Cover problem as a special case. When the distances form a metric, that is, they satisfy the triangle inequality, there are algorithms with constant approximation ratio. Shmoys, Tardos and Aardal [25] were the first ones to obtain an $O(1)$ -approximation algorithm for UFL. Currently the best known approximation ratio is 1.488 by Li [22]. On the other hand, Guha and Khuller [13] showed a lower bound of 1.463 on approximability, assuming $\text{NP} \not\subseteq \text{DTIME}(n^{O(\log \log n)})$. Sviridenko [28] strengthened the underlying assumption to $\text{P} \neq \text{NP}$.

The FTFP Problem The problem studied in this thesis is the Fault-Tolerant Facility Placement problem (FTFP), which generalizes the UFL problem. The difference from the UFL problem is that, in the FTFP problem, each client now has a demand that could be more than one, and each unit demand of a client needs to be connected to a different facility; at each site, we can open one or more facilities. Opening multiple facilities at the same site incurs a cost of the opening cost for this site multiplied by the number of facilities opened. The connection cost between a site and a client is the number of connections times

Shmoys, Tardos and Aardal	LP-rounding	3.16	1997 [25]
Chudak and Shmoys	LP-rounding	1.736	1998 [10]
Sviridenko	LP-rounding	1.582	2002 [26]
Jain and Vazirani	primal-dual	3	2001 [20]
Jain <i>et al.</i>	dual-fitting	1.61	2003 [18]
Arya <i>et al.</i>	local-search	3	2001 [1]
Byrka and Aardal	hybrid	1.50	2007 [4]
Li	hybrid	1.488	2011 (best result) [22]
Guha and Khuller	lower-bound	1.463	1998 [13]

Table 1.1: A history of approximation algorithms and their ratio for the UFL problem. The columns are author, technique, ratio and year. The word *hybrid* means the use of two independent algorithms and return the better solution of the two. The last row of 1.463 is about inapproximability, or hardness.

the distance between the two, with the constraint that the number of connections cannot exceed the number of open facilities at that site. The FTFP problem asks for a solution with minimum total cost, that is, the sum of facility cost and connection cost.

Like the UFL problem, the FTFP problem is NP-hard, and thus precludes a promise for polynomial-time algorithms that solve the problem exactly. Nonetheless, there are polynomial-time algorithms that deliver a solution whose cost is only a small percentage off from the optimal solution's cost. These algorithms, known as approximation algorithms, and their performance analysis, are the subject of this thesis.

1.2 Hardness Results on Approximating UFL and FTFP

Not only is the FTFP problem NP-hard, more can be said about the extent to which we can approximate an optimal solution. Results showing that problems cannot be approximated to certain numbers, using well repected assumptions like $P \neq NP$, are called hardness results and we present the hardness results on the FTFP problem in this section.

Since the fault-tolerant facility placement problem (FTFP) contains the uncapacitated facility location problem (UFL) as a special case (setting all $r_j = 1$ in FTFP we get UFL), any hardness result obtained on UFL remains applicable to FTFP. In the following, we review some well-known hardness results on UFL, with the implication that the same claims hold for FTFP as well.

Theorem 1 *The general¹ UFL problem is NP-hard.*

¹*general* means no restriction on the distances d_{ij} . This is in contrast with the *metric* version where d_{ij} satisfies the triangle inequality.

Proof. The proof is by reduction from the Set Cover problem. In the Set Cover problem, we are given a universe of elements, that is $\mathcal{U} = \{e_1, \dots, e_n\}$, and a collection of sets $\mathcal{S} = \{S_1, \dots, S_m\}$ such that $S_i \subseteq \mathcal{U}$ for $i = 1, \dots, m$. The problem asks for a set cover with minimum size, where a set cover is a collection of sets from \mathcal{S} whose union is \mathcal{U} . We construct a general UFL instance like this: for each element $e_j, j = 1, \dots, n$, we have a client j ; and for each set $S_i, i = 1, \dots, m$, we have a facility i . The facility cost $f_i = 1$ for every facility $i = 1, \dots, m$ ². The distance $d_{ij} = 1$ if $e_j \in S_i$ and $d_{ij} = \infty$ if $e_j \notin S_i$. Clearly an optimal solution for the UFL instance can only use connections with $d_{ij} = 1$. It is easy to see that given any optimal solution to the Set Cover instance, we can construct an optimal solution for the UFL instance, simply by taking the facilities whose corresponding sets are chosen in the set cover. On the other hand, given an optimal solution to the UFL instance, we can only have $d_{ij} = 1$ connections. This implies that, for every client j , suppose it was connected to facility i in the UFL solution, then the corresponding element e_j is covered by the corresponding set S_i . Let I be the set of facilities chosen to open in the UFL solution; it is easily seen that the corresponding sets in the Set Cover instance form a set cover. ■

Theorem 2 *The metric UFL problem is NP-hard.*

Proof. The reduction is also from the Set Cover problem. Suppose the Set Cover instance has universe $\mathcal{U} = \{e_1, \dots, e_n\}$ and collection of sets $\mathcal{S} = \{S_1, \dots, S_m\}$. Unlike the general UFL problem, we can no longer have edges with length 1 and ∞ now, as the distances are constrained by the triangle inequality. Our facility cost $f_i = \epsilon = 1/m^2$, and every facility i has the same facility cost. Our distance d_{ij} is now 1 if $e_j \in S_i$ and 3 if $e_j \notin S_i$.

²Actually any value of $f_i > 0$ will work, for example, we can set $f_i = 100$ for every facility i .

Given the construction, it is clear that any optimal solution for the metric UFL instance cannot use a connection of distance 3, as there exists another solution that beats such a solution with lower cost, namely a solution that opens all facilities with total cost $m\epsilon + n \cdot 1 = m\epsilon + n$, as we have chosen $\epsilon = 1/m^2$. It follows that any optimal solution for the metric UFL instance must have all clients connected at distance 1. Such a solution would then have a facility set corresponds to a set cover for the Set Cover instance. ■

Now we show the **MaxSNP**-hardness of the metric UFL problem. This implies that, under the assumption that $P \neq NP$, there exists some constant c such that the metric UFL problem cannot be approximated to better than c -approximation. As a consequence, the metric UFL problem cannot have polynomial-time approximation scheme (PTAS). PTAS are algorithms that, for any constant $\epsilon > 0$, compute a solution with cost no more than $(1 + \epsilon)$ from the optimal. Moreover, the running time is polynomial in the input size, with ϵ treated as a constant. The following theorem is due to Guha and Khuller [13].

Theorem 3 [13] *The metric UFL problem is MaxSNP-hard.*

Proof. The full proof can be found in [13] and we only sketch the main idea here. The proof is by a reduction from the B-Vertex Cover problem, a problem known to be **MaxSNP**-hard. In the B-Vertex Cover problem, we are given a graph $G = (V, E)$, and a constant B , such that every vertex $u \in V$ has degree no more than B . And the problem asks for a vertex cover with minimum size. That is, we are to find a minimum set $V' \subseteq V$ such that every edge $e \in E$ has at least one endpoint in V' .

The idea is to show that, for any given constant $0 < \epsilon < 1$, given an algorithm for the metric UFL problem with approximation ratio $1 + \epsilon$, we are able to find an algorithm

for the B-Vertex Cover problem with approximation ratio $1 + \epsilon'$ such that ϵ' is a constant depending on ϵ and possibly B , and ϵ' approaches 0 as ϵ approaches 0. It turns out that we can set $\epsilon' = (1 + B)\epsilon$ for our purpose. ■

From Theorem 3, we know there exists a constant c such that the UFL problem cannot be approximated to better than c , assuming $P \neq NP$. The last piece of hardness result presented in this section gives the best known such constant $c = 1.463$. This result is also due to Guha and Khuller [13], with an improvement from Sviridenko [28].

Theorem 4 [13, 28] *UFL cannot be approximated to less than 1.463 unless $P = NP$.*

Proof. The proof is by contradiction. More precisely, we show that if the metric UFL problem can be solved by a polynomial-time algorithm with approximation ratio γ that is less than $\gamma_0 = 1.463$ ³, then we would have a polynomial-time algorithm with approximation ratio $(1/\rho)\ln n$ for some constant $\rho > 1$ for the Set Cover problem. Here n is the number of elements in the universe in the Set Cover instance. Using a result by Feige [12], the existence of a $(1/\rho)\ln n$ approximation algorithm for a constant $\rho > 1$ implies $NP \subseteq DTIME(n^{O(\log \log n)})$.

Given a Set Cover instance with a universe $\mathcal{U} = \{e_1, \dots, e_n\}$ of elements and a family of sets $\mathcal{S} = \{S_1, \dots, S_m\}$ with every set $S_i \in \mathcal{S}$ being a subset of \mathcal{U} . The proof proceeds in iterations. In each iteration we construct a metric UFL instance and run the supposedly γ -approximation algorithm for UFL. After that we remove clients connected at distance of 1 (call them *covered* in this iteration) in the UFL solution. So future iterations will deal with a smaller set of clients. We repeat until all clients are covered.

³ γ_0 is the solution to the equation $\gamma = 1 + 2/e^\gamma$.

Consider the t^{th} iteration and let n_t be the number of clients not yet covered. The metric UFL instance for this iteration consists of a set of facilities \mathbb{F} corresponding to the set \mathcal{S} , and a set of clients with one client j for each uncovered element e_j . The distance d_{ij} is 1 if $e_j \in S_i$, and 3 otherwise. The opening cost of all facilities are equal and we set every $f_i = cn_t/k$, where c is some constant decided later to optimize our lower bound on approximation ratio. Here k is the number of sets in an optimal solution for the Set Cover instance. Note that we can perform the rest of the proof for every possible $k = 1, \dots, n$, so we can assume the knowledge of k .

Our construction ensures that, at any iteration, if the UFL solution found by the γ -approximation UFL algorithm does not cover a large portion of the clients, then the ratio γ between this UFL solution's cost and the optimal solution's cost must be at least 1.463. The solution with k facilities and all clients connected at distance 1 is used for estimating the optimal UFL solution's cost. The other case is that in every iteration we have an UFL solution that covers a large portion of clients. This would then give us a solution for the given Set Cover instance with cost (number of sets) no more than $(k/\rho) \ln n$, that is no more than $(1/\rho) \ln n$ times of an optimal solution for constant $\rho > 1$. This then implies $\text{NP} \subseteq \text{DTIME}(n^{O(\log \log n)})$ by Feige's result [12].

Using an observation by Sviridenko, the underlying assumption $\text{NP} \not\subseteq \text{DTIME}(n^{O(\log \log n)})$ can be strengthened to $\text{P} \neq \text{NP}$. That is, metric UFL cannot have a polynomial-time algorithm with ratio less than 1.463 unless $\text{P} = \text{NP}$. ■

With that, we conclude the discussion on the hardness results for the Uncapaciated Facility Location problem (UFL), with the implication that all these hardness results carry

on to our problem, the Fault-tolerant Facility Placement problem (FTFP) as well.

Chapter 2

Related Work and Known Results

In this chapter we review the history and results of the Fault-Tolerant Facility Placement problem (FTFP), and those for the two closely related problems, the Uncapacitated Facility Location problem (UFL) and the Fault-Tolerant Facility Location problem (FTFL). In all three problems, we are given a set of sites \mathbb{F} that we could open facilities, and a set of clients \mathbb{C} with demands that need to be connected to open facilities. The cost to open one facility at site i is f_i , and the cost to make one connection between a client j to a facility at site i is d_{ij} . A client j with demand r_j needs to be connected to r_j different facilities (Facilities at the same site are considered different.). Our goal is to open a set of facilities at the given sites, and connecting clients to those facilities in such a way that all clients' demands are satisfied, and the total cost is minimized. The difference between the problems are: in the UFL problem, all demands are 1, that is $r_j = 1$ for all clients j , then we need no more than 1 facility at a site; In the FTFL problem, demands may be more than 1, but each site can open at most 1 facility; In the FTFP problem, demands may be

more than 1, and each site can open any number of facilities, including zero.

In designing algorithms for the facility location problems (FL), which encompass all three problems, UFL, FTFL, and FTFP, we have two competing goals: on one hand we want to open as few facilities as possible so that our facility cost is small; On the other hand we need as many facilities as possible so that every client can connect to nearby facilities.

Before delving into specific results and techniques for the facility location problems, we give some simple observations about the UFL problem, which apply to other FL problems with minor adjustments. In the UFL we need no more than 1 facility in each site so we talk about facility without mentioning the site. The following observations are obvious.

Observation 5 *If all facility opening cost are zero, then an optimal solution is to open all facilities and connect each client to the nearest facility.*

Observation 6 *If all distances are zero, that is, all facilities and clients are colocated at the same point, then an optimal solution is to open the cheapest facility and connect all clients to that facility.*

Observation 7 *For a fixed set of open facilities, there is a polynomial algorithm to find the optimal way to connect clients to facilities.*

Proof. For UFL, it is a simple matter of finding the nearest open facility for each client and a brute force search only takes time $O(|F| \cdot |C|)$. Even for more complicated variants where the facilities each has a capacity, the optimal assignment can be found by using the mincost flow algorithm. ■

Next we review the known algorithms for UFL and FTFL. These two are well-studied problems in literature, and are closely related to our problem, FTFP. In particular, the LP-rounding algorithms for UFL inspired our approach to the FTFP problem.

2.1 Related Work on UFL

The Uncapacitated Facility Location problem (UFL) is the simplest variant of the Facility Location problems (FL), and has received the most attention in research community. One thing that is surprising about UFL is that a wide range of different techniques to design approximation algorithms have been successful on the UFL problem. As is the case for many other optimization problems, the analysis of the approximation ratio hinges on estimating the optimal value using the solution to a linear program (LP). For the UFL problem, the LP formulation by Balinski [2] is now standard. In the LP we use $y_i \in \{0, 1\}$ to indicate whether facility $i \in \mathbb{F}$ is open or not, and $x_{ij} \in \{0, 1\}$ to indicate whether client j is connected to facility i . Relaxing the integral constraint to allow fractional values, we obtain the following LP (2.1) for the UFL problem. Observe that we do not need an explicit constraint of $x_{ij} \leq 1$ or $y_i \leq 1$ as any optimal solution to LP (2.1) must satisfy the two constraints automatically.

The LP is,

$$\begin{aligned}
& \text{minimize} && \sum_{i \in \mathbb{F}} f_i y_i + \sum_{i \in \mathbb{F}, j \in \mathbb{C}} d_{ij} x_{ij} && (2.1) \\
& \text{subject to} && y_i - x_{ij} \geq 0 && \forall i \in \mathbb{F}, j \in \mathbb{C} \\
& && \sum_{i \in \mathbb{F}} x_{ij} \geq 1 && \forall j \in \mathbb{C} \\
& && x_{ij} \geq 0, y_i \geq 0 && \forall i \in \mathbb{F}, j \in \mathbb{C}
\end{aligned}$$

The dual program is,

$$\begin{aligned}
& \text{maximize} && \sum_{j \in \mathbb{C}} \alpha_j && (2.2) \\
& \text{subject to} && \sum_{j \in \mathbb{C}} \beta_{ij} \leq f_i && \forall i \in \mathbb{F} \\
& && \alpha_j - \beta_{ij} \leq d_{ij} && \forall i \in \mathbb{F}, j \in \mathbb{C} \\
& && \alpha_j \geq 0, \beta_{ij} \geq 0 && \forall i \in \mathbb{F}, j \in \mathbb{C}
\end{aligned}$$

Two general approaches using LP to design approximation algorithms are LP-rounding and primal-dual. LP-rounding algorithms start with calling an LP-solver to obtain an optimal fractional solution $(\mathbf{x}^*, \mathbf{y}^*)$, and then round the fractional solution in a way that preserves the feasibility while keeps the cost from increasing by too much. On the other hand, primal-dual algorithms do not require solving the LP and the use of LP in the algorithms and their analysis is implicit. Those algorithms work by constructing an

integral feasible primal solution and a feasible (fractional) dual solution simultaneously, such that the cost of the primal solution is related to the cost of the dual solution. We can then estimate the cost of an optimal solution by the cost of the dual solution. Note that the cost of any feasible dual solution provides a lower bound of the optimal fractional solution value for the primal, assuming the primal is a minimization program and the dual is a maximization program. For more information about Linear Programming and its application in approximation algorithms, we refer the reader to Appendix A.1.

Approximation Results The first $O(1)$ -approximation algorithm was obtained by Shmoys, Tardos and Aardal [25], with a ratio of 3.16, using LP-rounding. The Shmoys *et al.*'s algorithm also established a general framework that underpins all subsequent LP-rounding algorithms. In the framework the clients are partitioning into clusters and each cluster has a representative client. The rounding algorithm guarantees that each representative has a nearby facility to connect to, and the rest of the clients can then use the facility via their representative. The same solution structure is used in all known LP-rounding algorithms for the UFL problem. The 3.16 algorithm had made some greedy choices in several steps and left quite some room for improvement. Chudak and Shmoys [9] was the first to use the idea of randomized rounding to improve the approximation ratio. Roughly speaking, in their algorithm, each facility i is opened with probability y_i^* where y_i^* is given by the optimal fractional solution $(\mathbf{x}^*, \mathbf{y}^*)$ to the LP (2.1). The expected connection cost is then estimated using a provably worse random process in which each facility was opened independently and hence the expected connection cost is easier to analyze. With a very refined analysis and some technical inequality, they obtained a ratio of $1 + 2/e = 1.736$. Sviridenko [26] further

improved the ratio to 1.582, by using a scaled version of the optimal fractional solution, and a judicious choice of some distribution of the scaling parameter. The rounding process is called pipage rounding, a deterministic rounding process that takes advantage of the concave property of some cost function. The analysis is highly technical.

Combinatorial algorithms that do not require solving the LP have also been found successful in obtaining good approximation ratio for UFL. Most notably, Jain and Vazirani [20] introduced a primal-dual algorithm with ratio 3. The algorithm grows a feasible dual solution from zero, and updates a primal solution accordingly until the primal solution is feasible. The approximation ratio is obtained via a relaxed version of the complementary slackness conditions. The conditions provide a bound on the cost of the primal solution in terms of the cost of the dual solution. Clearly the cost of a feasible dual solution provides a lower bound on the cost of an optimal primal solution, assuming the primal is a minimization problem. A slightly different primal-dual based result is obtained by Jain, Markakis, Mahdian, Saberi and Vazirani [18]. They analyzed a greedy algorithm that repeatedly picks the most cost-effective star until all clients are connected. A star consists of one facility and a subset of clients. The cost-effectiveness is the cost of the star, which includes facility cost and connection cost of all clients to that facility, divided by the number of clients in that star. In each iteration, the algorithm picks the best star, connects all member clients to the facility, and set the facility cost to zero. The clients just get connected are removed from future consideration but the facility could be reused for future stars. Jain *et al.* analyzed the greedy algorithm and its variant using the dual-fitting technique. They first showed that the greedy algorithm can be interpreted as a process of growing a dual solution and

updating a primal solution. Moreover, the cost of the primal solution is equal to the cost of the dual solution produced by the greedy algorithm. It might appear that we have solved the UFL problem optimally, although we know that this cannot be the case, as the UFL problem is NP-hard. The catch is that the dual solution computed by the greedy algorithm is not feasible. The next step is to find a common factor γ , such that the dual values, after shrinking by γ , is feasible to the dual program. That common factor γ is then the desired approximation ratio. They showed their algorithms have approximation ratio 1.861 ¹ and 1.61 respectively.

A still different approach is local search, in which we start with a feasible integral solution and make local moves to improve the solution, and stop at some local optimum. The set of allowed local moves needs to be chosen carefully, as more powerful moves allow local optimal to be closer to global optimal, on the other hand restricting on a few simple moves results in faster algorithms and easier analysis. Arya *et al.* [1] showed that local search gives a ratio of 3 for the UFL problem.

To date, the best known approximation results for UFL are due to Byrka [4] with ratio 1.5 and a follow-up work by Li [22] with ratio 1.488. Both use a combination of two algorithms, one is an LP-rounding algorithm and the other is the 1.61 greedy algorithm of Jain, Mahdian and Saberi [19]. Since the hybrid approach requires the introduction of the notion of *bifactor* analysis, we postpone the discussion of Byrka and Li's result for now and introduce the notion of bifactor analysis first.

Bifactor Analysis Given that the total cost of a solution to the UFL problem consists of

¹An improved analysis by Mahdian showed the actual ratio is 1.81.

two parts, the facility cost and the connection cost, a notion of bifactor approximation was introduced by Jain *et al.* in [18]. An algorithm with facility cost F_{ALG} (sum of f_i for all open facility i) and connection cost C_{ALG} (sum of d_{ij} for pairs of (i, j) connected), is said to be (γ_f, γ_c) -approximation if, for every feasible solution SOL, with facility cost F_{SOL} and connection cost C_{SOL} , we have

$$F_{\text{ALG}} + C_{\text{ALG}} \leq \gamma_f F_{\text{SOL}} + \gamma_c C_{\text{SOL}}.$$

In particular, for the optimal fractional solution $(\mathbf{x}^*, \mathbf{y}^*)$ with facility cost $F^* = \sum_{i \in \mathbb{F}} f_i y_i^*$ and connection cost $C^* = \sum_{j \in \mathbb{C}} d_{ij} x_{ij}^*$, we have

$$F_{\text{ALG}} + C_{\text{ALG}} \leq \gamma_f F^* + \gamma_c C^*.$$

The notion of bifactor approximation is helpful when an algorithm has imbalanced factors γ_f and γ_c . It is easy to see such an algorithm has approximation ratio $\max\{\gamma_f, \gamma_c\}$. However, more can be said, as there are techniques like cost scaling and greedy augmentation to balance the two factors, that is, to decrease one at the expense of increasing the other, thus achieving a better overall approximation ratio. The techniques of cost scaling and greedy augmentation and their use to balance the two factors are introduced by Guha, Khuller and Charikar [13, 8]. For example, the primal-dual algorithm by Jain and Vazirani [21] is a (1,3)-approximation algorithm and using cost scaling and greedy augmentation, it is possible to show the algorithm can achieve a ratio of 1.85 [8].

Finally we return to the hybrid algorithms by Byrka and Aardal [5] and Li [22]. Byrka and Aardal gave an LP-rounding algorithm with bifactor (1.68, 1.37), and showed that this algorithm, when combined with a (1.11, 1.78) algorithm by Jain *et al.* [18], gives

a ratio of 1.50. Li showed that by choosing a nontrivial distribution of the scaling factor of Byrka's algorithm, the analysis can be refined to show an overall ratio of 1.488. The 1.488 ratio is currently best known approximation result.

LP-rounding on UFL We now present a more detailed description on the LP-rounding approaches, as our results on FTFP are built on the work of LP-rounding for UFL.

Every LP-rounding algorithm for UFL starts with solving the LP (2.1) to obtain an optimal fraction solution $(\mathbf{x}^*, \mathbf{y}^*)$. Then we need to round the fractional solution to an integral solution $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$ without increasing the cost by much. As hinted in the above discussion, an integral solution with small cost would have each client connect to a nearby facility and few facilities open. Consider a client j , to get a handle on the connection cost, we would like j to connect to some neighboring facility $i \in N(j)$, where $N(j) \stackrel{\text{def}}{=} \{i \in \mathbb{F} : x_{ij}^* > 0\}$. For the sake of connection cost, it is desirable for a client to have a neighboring facility open, as those are facilities not too far away. In fact, at the end of Appendix A.1 we give a bound on the distance from a client to any of its neighbors. However, it is in general not possible to have every client connect to a neighboring facility, or we would have to open too many facilities, and thus incur a high facility cost. An alternative is to select a subset of clients, denoted as $C' \subseteq \mathbb{C}$ and only require clients in C' have a neighboring facility open. Clients outside C' are then connected to a facility via some client in C' . The connection cost for clients in $\mathbb{C} \setminus C'$ are then bounded using the triangle inequality. For this strategy to work, the clients j outside C' need to be able to find some client j' in C' such that both $d_{jj'} \stackrel{\text{def}}{=} \min_{i \in \mathbb{F}} d_{ij} + d_{ij'}$ and $d_{\phi(j')j'}$ are small. Here $\phi(j')$ is the facility that j' connects to.

Suppose we have chosen a subset C' of client, we then open one facility for each

client j' in C' and j' would connect to that facility. Let j be any client in \mathbb{C} , denote by $\phi(j)$ the facility that j connects to. To bound the facility cost F_{ALG} , it suffices to require the clients j' in C' have disjoint neighborhood $N(j')$, and $f_{\phi(j')}$ be bounded by the average facility cost of $N(j')$, that is, $f_{\phi(j')} \leq \sum_{i \in N(j')} f_i y_i^*$. Call two clients j_1 and j_2 *related* if their neighborhood overlap, that is $N(j_1) \cap N(j_2) \neq \emptyset$. The requirement of disjoint neighborhood on the set C' immediately implies that between every two related clients, we can select at most one of them into the set C' . Suppose we choose a client j' into the set C' . This would then disqualify every client j that is related to j' . We then define a cluster as the set of all clients j that are related to j' , and call j' the *representative* of that cluster, and we say other clients in that cluster are *assigned* to that representative.

Suppose we have chosen the set C' of representative clients, and we open one facility, which is $\phi(j')$ for each representative client j' in C' , in a way that $f_{\phi(j')}$ is no more than the average cost of facilities in $N(j')$, that is $f_{\phi(j')} \leq \sum_{i \in N(j')} f_i y_i$. Then we know that the facility cost is bounded by $\sum_{i \in \mathbb{F}} f_i y_i^* = F^*$ because the neighborhood of clients in C' are disjoint.

For connection cost, each client j' in C' connects to the only facility open in $N(j')$ and we have mentioned earlier that there is a way to bound the distance from $\phi(j')$ to j' when $\phi(j') \in N(j')$, see the end of Appendix A.1. For clients j not in C' , we need to have an open facility for it to connect to. One possibility is to use $\phi(j')$, if j is assigned to j' . When that happens, $d_{\phi(j'),j}$ can be bounded by the triangle inequality as

$$d_{\phi(j'),j} = d_{\phi(j'),j} \leq d_{\phi(j'),j'} + d_{j',j} \leq d_{\phi(j'),j'} + d_{i,j'} + d_{i,j} \text{ for all } i \in \mathbb{F}. \quad (2.3)$$

In the above $d_{jj'}$ denotes the shortest distance between client j and j' , that is $d_{jj'} =$

$\min_{i \in \mathbb{F}} d_{ij'} + d_{ij}$. Now if we can show a bound on both $d_{\phi(j'),j'}$ and $d_{j'j}$ in terms of some quantity of client j , we can expect to have a bound on $d_{\phi(j),j}$ using the same quantity. Since the inequality (2.3) holds for any facility $i \in \mathbb{F}$, it would certainly hold for some facility i' in $N(j') \cap N(j)$. Such facility i' exists because of the way the relation between clients are defined. We then have a bound on all three quantities, namely $d_{\phi(j'),j'}$, $d_{i'j'}$ and $d_{i'j}$, because the corresponding facilities are all in the clients' neighborhood. Now we only need $d_{\phi(j')j'}$ and $d_{i'j'}$ to be small in terms of $d_{i'j}$. Since we have a choice on which representatives to pick from the client set, we greedily pick a representative client to minimize some measure on a client, and then remove both the client and clients assigned to it from future consideration. Repeat until no client remains to be considered. This process gives us a set of representatives and their associated clusters of clients. This structure is then used in the rounding process.

We now describe a general rounding process based on the above discussions. Let $N'(j)$ be a subset of facilities defined for each client j and let $g(j)$ be some measure on client j such that for any facility i in $N'(j)$, we have $d_{ij} \leq g(j)$. For any reasonable choice of $g(j)$, we shall have $N'(j)$ be a subset of $N(j)$, although $N'(j)$ could be a strict subset. Recall that $N(j) \stackrel{\text{def}}{=} \{i \in \mathbb{F} : x_{ij}^* > 0\}$ is the neighborhood of client j . Then our integral solution consists of one facility from each reduced neighborhood $N'(j')$. Notice that this implies that every representative client j' is connected to a facility with distance at most $g(j')$. Given $g(j)$ for all clients $j \in \mathbb{C}$, we can choose the representative j' as the one that minimizes the measure $g(j)$ among all clients assigned to j' ². One consequence is that any client j outside C' , that is, clients that are not representatives, can now have their connection cost

²We consider j' assigned to itself.

being bounded as:

$$d_{\phi(j),j} \leq d_{\phi(j'),j'} + d_{i'j'} + d_{ij} \leq g(j') + g(j') + g(j) \leq g(j) + g(j) + g(j) = 3g(j).$$

In other words, given the function $g(j)$ and its associated reduced neighborhood $N'(j)$ for each client $j \in \mathbb{C}$, we are able to round the fractional solution $(\mathbf{x}^*, \mathbf{y}^*)$ into an integral solution such that $F_{\text{ALG}} \leq F^*$ and $d_{\phi(j),j} \leq 3g(j)$. Summing the second inequality over all clients, we have the total connection cost bounded by $3 \sum_{j \in \mathbb{C}} g(j)$.

Next we discuss specific choices of the function $g(j)$. Two ways are possible: one is to use the dual solution $(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*)$, and then the complementary slackness conditions tell us that $d_{ij} \leq \alpha_j^*$ for every $i \in N(j) = \{i \in \mathbb{F} : x_{ij}^* > 0\}$. This allows us to define $g(j) = \alpha_j^*$ and $N'(j) = N(j)$. The other way is to trim the neighborhood of a client by throwing away most distant neighbors so that the farthest facility remaining in the trimmed neighborhood is at a distance within some factor of the average distance $C_j^* \stackrel{\text{def}}{=} \sum_{i \in N(j)} d_{ij} x_{ij}^*$. We now open one facility in $N'(j')$ for each representative $j' \in C'$ in a way that $f_{\phi(j')}$ is no more than the average facility cost of facilities in $N'(j')$, that is $\sum_{i \in N'(j')} f_i y_i / \sum_{i \in N'(j')} y_i$. Now we have a bound on $d_{\phi(j),j}$ in terms of C_j^* . However, our facility cost is now bounded by

$$F_{\text{ALG}} \leq \sum_{j' \in C'} \left(\sum_{i \in N'(j')} f_i y_i^* / \sum_{i \in N'(j')} y_i^* \right).$$

To estimate this cost, we assume a lower bound γ on $\sum_{i \in N'(j')} y_i^*$ for every client $j' \in C'$.

Our facility cost can then be bounded as

$$\begin{aligned}
F_{\text{ALG}} &\leq \sum_{j' \in C'} \left(\sum_{i \in N'(j')} f_i y_i^* / \sum_{i \in N'(j')} y_i^* \right) \quad (\text{choice of } \phi(j')) \\
&\leq \sum_{j' \in C'} \sum_{i \in N'(j')} f_i y_i^* / \gamma \quad (\text{assumption } \sum_{i \in N'(j')} y_i^* \leq \gamma) \\
&= \frac{1}{\gamma} \sum_{j' \in C'} \sum_{i \in N'(j')} f_i y_i^* \\
&\leq \frac{1}{\gamma} \sum_{j' \in C'} \sum_{i \in N(j')} f_i y_i^* \quad (\text{replace } N'(j) \text{ with } N(j)) \\
&= \frac{1}{\gamma} \sum_{i \in \mathbb{F}} f_i y_i^* = F^*.
\end{aligned}$$

We note that the factor that bounds $g(j)$ by C_j^* is related to γ and we denote that factor as a function of γ , that is $f(\gamma)$. Thus we have $g(j) \leq f(\gamma)C_j^*$. It follows that we can obtain an integral solution with facility cost no more than γF^* and connection cost no more than $3 \sum_{j \in \mathbb{C}} g(j) = 3 \sum_{j \in \mathbb{C}} f(\gamma)C_j^* = 3f(\gamma)C^*$. Recall that $F^* = \sum_{i \in \mathbb{F}} f_i y_i^*$ and $C^* = \sum_{i \in \mathbb{F}, j \in \mathbb{C}} d_{ij} x_{ij}^*$ and we also define $\text{LP}^* = F^* + C^*$. Since $\text{LP}^* \leq \text{OPT}$, because the cost of an optimal fractional solution, that is LP^* is no more than that of an optimal integral solution OPT . Our rounded solution then has approximation ratio $\max\{\gamma, 3f(\gamma)\}$.

For concreteness we now describe a simple rounding algorithm by Chudak [10], as well as the original algorithm by Shmoys, Tardos and Aardal [25], followed by improvement by Chudak and Shmoys [9] and Sviridenko [26], Byrka [4] and Li [22].

We start with an algorithm by Chudak that gives an easy 4-approximation. The algorithm starts with fractional optimal solution $(\mathbf{x}^*, \mathbf{y}^*)$ for the primal and optimal solution

(α^*, β^*) for the dual. The algorithm then repeatedly picks an unconnected client j with minimum α_j^* , open a cheapest facility i in $N(j)$, and then connect j to i . In addition, unconnected clients with a neighborhood overlapping $N(j)$ are connected to i as well. The algorithm terminates when all clients are connected. Given the previous discussion, clearly the facility cost is no more than F^* . For connection cost, we can use $g(j) = \alpha_j^*$, and every client j is connected to a facility with distance no more than $3g(j) = 3\alpha_j^*$. Summarizing, we have total cost no more than $F^* + \sum_{j \in \mathbb{C}} 3\alpha_j^* = F^* + 3\text{LP}^* \leq 4\text{LP}^* \leq 4\text{OPT}$, where LP^* is the cost of an optimal fractional solution and OPT is the cost of an optimal integral solution.

Another 4-approximation algorithm, the original approach by Shmoys, Tardos and Aardal, uses a different way to obtain $g(j)$, by cutting the neighborhood $N(j)$ to a smaller set $N'(j)$. More precisely, for each client j , they cut at a distance $g(j)$ such that the accumulated fractional value is at least some constant γ , that is

$$\sum_{i \in N(j) : d_{ij} \leq g(j)} x_{ij}^* \geq \gamma.$$

By definition, the farthest facility in the trimmed neighborhood $N'(j)$, is at a distance of $g(j)$ from client j . To bound $g(j)$, we use an argument similar to the proof of Markov's inequality. More precisely, we have

$$C_j^* \stackrel{\text{def}}{=} \sum_{i \in N(j)} d_{ij} x_{ij}^* \geq \sum_{i \in N(j) \setminus N'(j)} d_{ij} x_{ij}^* \geq \sum_{i \in N(j) \setminus N'(j)} g(j) x_{ij}^* \geq g(j)(1 - \gamma).$$

As a result $g(j) \leq C_j^*/(1 - \gamma)$ and notice that $\sum_{j \in \mathbb{C}} C_j^* = C^*$. As discussed above, we have total connection cost no more than $3C^*/(1 - \gamma)$. For the facility cost, since we have limited ourselves only facilities in $N'(j)$, as opposed to $N(j)$. If we open the cheapest facility in

$N'(j)$ for each chosen client j , all we can say now is that the facility cost is not more than

$$\sum_{j \in P} \sum_{i \in N'(j)} f_i y_i^* / \sum_{i \in N'(j)} y_i^*.$$

Here P is the set of representative clients that were chosen by the algorithm. We are blessed since we have an lower bound of γ on $\sum_{i \in N'(j)} y_i^*$ for every client j . As shown earlier, the above sum can be upper bounded by

$$\sum_{j \in P} \sum_{i \in N'(j)} f_i y_i^* / \gamma \leq \frac{1}{\gamma} \sum_{j \in P} \sum_{i \in N'(j)} f_i y_i^* \leq \frac{1}{\gamma} \sum_{i \in \mathbb{F}} f_i y_i^* = \frac{1}{\gamma} F^*.$$

So our approximation ratio is $\max\{3/(1 - \gamma), 1/\gamma\}$. Pick $\gamma = 1/4$ and we have a 4-approximation.

The simplest rounding algorithms described so far is clearly not the best possible and there have been several improvement proposed afterwards. Nonetheless, they all follow the same structure by partitioning clients into clusters with representatives, and open exactly one facility for each representative client. To save on connection cost, one observes that a not chosen client can use a neighboring facility provided that the neighbor is open. For facility cost, it is not necessary to open the cheapest facility in the reduced neighborhood $N'(j)$. By choosing each facility in a neighborhood with probability proportional to y_i^* , we would have an expectation of facility cost bounded by F^* as well, and this would help the expected connection cost. The analysis, however, is much more involved and requires one to estimate the probability of a non-representative client connecting to its neighbor and the probability to connect to the facility opened by its assigned representative. We shall not elaborate the intricacies here as our rounding algorithms for the FTFP problem in the main part, Chapter 5 will rephrase and prove similar theorems there. This concludes our

review on the known approximation results for the Uncapacitated Facility Location problem (UFL).

2.2 Related Work on FTFL

The fault-tolerant facility location problem (FTFL), was first introduced by Jain and Vazirani [21]. The LP is

$$\begin{aligned}
& \text{minimize} && \sum_{i \in \mathbb{F}} f_i y_i + \sum_{i \in \mathbb{F}, j \in \mathbb{C}} d_{ij} x_{ij} && (2.4) \\
& \text{subject to} && y_i - x_{ij} \geq 0 && \forall i \in \mathbb{F}, j \in \mathbb{C} \\
& && \sum_{i \in \mathbb{F}} x_{ij} \geq r_j && \forall j \in \mathbb{C} \\
& && y_i \leq 1 && \forall i \in \mathbb{F} \\
& && x_{ij} \geq 0, y_i \geq 0 && \forall i \in \mathbb{F}, j \in \mathbb{C}
\end{aligned}$$

Jain and Vazirani adapted their primal-dual algorithm for UFL to FTFL to gain a ratio of $3 \ln R$ where $R = \max_j r_j$ is the maximum demand among all clients. The first constant approximation algorithm was given by Guha, Meyerson and Munagala [14], using LP-rounding similar to the Shmoys, Tardos and Aardal's [25] approach for the UFL problem. Subsequent improvement was made by Swamy and Shmoys [27] using pipage rounding with a ratio of 2.076. The current best known approximation ratio is 1.7245, due to Byrka, Srinivasan and Swamy [7], using dependent rounding with a laminar clustering structure.

We note that all the known $O(1)$ -approximation algorithms for FTFL are LP-

rounding algorithms and they need to solve the LP as a first step. Given the success of primal-dual based approaches for UFL, it is natural to ask whether such algorithms could be adapted to FTFL with a good ratio. To the best of the author’s knowledge, it is not known whether there is a primal-dual algorithm for FTFL with a sub-logarithmic approximation ratio. This is in stark contrast with the fact that two different primal-dual algorithms [21, 18] achieve constant ratio for UFL. We shall have more to say about this in Chapter 6.

2.3 Our Problem FTFP

Our problem, the fault-tolerant facility placement problem (FTFP), was introduced by Xu and Shen [29]³. The study of FTFP was partly motivated to obtain a better understanding of the implication of the fault-tolerant requirement on facility location problems. The Xu and Shen’s results seem to be valid only for a special case of FTFP. We later adapted the Chudak’s 4-approximation algorithm for UFL to FTFP, thus obtaining the first $O(1)$ -approximation algorithm for FTFP [30]. The algorithm was based on LP-rounding. Since then we have worked on improving the algorithm and analysis on LP-rounding. Our most recent results in [31] demonstrated that LP-rounding algorithms for FTFP can achieve an approximation ratio that matches that for UFL. For the applicability of primal-dual techniques on FTFP, we provide an explanation of possible difficulty in obtaining sub-logarithmic ratio using the dual-fitting approach. More details on primal-dual approaches can be found in Chapter 6.

³In their paper they call the problem the fault-tolerant facility allocation problem, or FTFA.

Chapter 3

Linear Program

3.1 Notation and Definition

In the fault-tolerant facility placement problem (FTFP), we denote the set of sites as \mathbb{F} and the set of clients as \mathbb{C} . Each client $j \in \mathbb{C}$ has a demand r_j , meaning that client j needs to be connected to r_j different facilities. The distance between a site i and a client j is denoted as d_{ij} . To open one facility at site i incurs a cost of f_i . To make one connection from client j to a facility at site i incurs a cost of d_{ij} . An FTFP instance is then fully specified by \mathbb{F} , \mathbb{C} , r_j for every $j \in \mathbb{C}$, and d_{ij} for every $i \in \mathbb{F}, j \in \mathbb{C}$.

A solution to the FTFP problem is a vector of (\mathbf{x}, \mathbf{y}) such that $x_{ij} \in \{0, 1, 2, \dots\}$ denotes the number of connections between site i and client j , and $y_i \in \{0, 1, 2, \dots\}$ denotes the number of facilities opened at site i . A feasible solution is one such that $y_i \geq x_{ij}$ for every $i \in \mathbb{F}, j \in \mathbb{C}$ and $\sum_{i \in \mathbb{F}} x_{ij} = r_j$ for all clients $j \in \mathbb{C}$. We would like to have a feasible solution with minimum total cost, that is $\sum_{i \in \mathbb{F}} f_i y_i + \sum_{i \in \mathbb{F}, j \in \mathbb{C}} d_{ij} x_{ij}$. For a given solution

(\mathbf{x}, \mathbf{y}) , we call the first term $\sum_{i \in \mathbb{F}} f_i y_i$ the facility cost of the solution, and the second term $\sum_{i \in \mathbb{F}, j \in \mathbb{C}} d_{ij} x_{ij}$ the connection cost of the solution.

For both the uncapacitated facility location problem (UFL) and the fault-tolerant facility placement problem (FTFP), we have a metric version of the problem in which the distances satisfy the triangle inequality, that is, for any two facilities/sites i_1, i_2 and any two clients j_1, j_2 , we have

$$d_{i_1 j_2} \leq d_{i_1 j_1} + d_{i_2 j_1} + d_{i_2 j_2}.$$

Unless otherwise stated explicitly, we assume the problem under study is the metric version.

In all our algorithms, we use the standard linear program (3.1) and its dual (3.2), first formulated by Balinski [2], to develop our algorithms and estimate the optimal solution value. For readers unfamiliar with Linear Programming (LP) and Integer Programming (IP), and their application on approximation algorithms, we included a short introductory section in the appendix A.1.

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.

$$\begin{aligned}
& \text{minimize} && \text{cost}(\mathbf{x}, \mathbf{y}) = \sum_{i \in \mathbb{F}} f_i y_i + \sum_{i \in \mathbb{F}, j \in \mathbb{C}} d_{ij} x_{ij} && (3.1) \\
& \text{subject to} && y_i - x_{ij} \geq 0 && \forall i \in \mathbb{F}, j \in \mathbb{C} \\
& && \sum_{i \in \mathbb{F}} x_{ij} \geq r_j && \forall j \in \mathbb{C} \\
& && x_{ij} \geq 0, y_i \geq 0 && \forall i \in \mathbb{F}, j \in \mathbb{C}
\end{aligned}$$

The dual program is,

$$\begin{aligned}
& \text{maximize} && \sum_{j \in \mathbb{C}} r_j \alpha_j && (3.2) \\
& \text{subject to} && \sum_{j \in \mathbb{C}} \beta_{ij} \leq f_i && \forall i \in \mathbb{F} \\
& && \alpha_j - \beta_{ij} \leq d_{ij} && \forall i \in \mathbb{F}, j \in \mathbb{C} \\
& && \alpha_j \geq 0, \beta_{ij} \geq 0 && \forall i \in \mathbb{F}, j \in \mathbb{C}
\end{aligned}$$

In each of our algorithms we will fix some optimal solutions of the LPs (3.1) and (3.2) 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 (3.1) and (3.2). 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 OPT for the optimal integral solution of (3.1). OPT is the value we wish to approximate, but, since $\text{OPT} \geq \text{LP}^*$, we can instead use 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 [9], 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 LP^* , as well as OPT , 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. Those properties allow a smooth generalization of known LP-rounding algo-

rithms for UFL to be applicable to FTFP, and handles the fault-tolerant requirement of the problem without effort.

Chapter 4

Techniques

After obtaining an optimal fractional solution to LP (3.1) and (3.2), we employ two techniques to obtain approximation results on the FTFP problem. Our first technique, which we call *demand reduction*, allows us to restrict our attention to a restricted version of the FTFP problem, in which all demands r_j are not too large. This restriction then sets stage for the application of our next technique, *adaptive partition*, so that we obtain an FTFP instance with facilities at sites and unit demand points derived from clients. For this FTFP instance, each facility can be either open or close, and each unit demand point connects to one of the open facilities. We would like to point out that we still need to cater the fault-tolerant requirement, that is, unit demands originated from the same client must connect to different facilities. We shall see that our adaptive partitioning step takes care of the fault-tolerant requirement smoothly.

4.1 Demand Reduction

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 9.) 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 (3.1). 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 3, 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

$$\hat{y}_i \leftarrow \lfloor y_i^* \rfloor, \quad \hat{x}_{ij} \leftarrow \lfloor x_{ij}^* \rfloor \quad \text{and}$$

$$\dot{y}_i \leftarrow y_i^* - \lfloor y_i^* \rfloor, \quad \dot{x}_{ij} \leftarrow x_{ij}^* - \lfloor x_{ij}^* \rfloor$$

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 8 (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 (3.1) 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 . ■

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 9 *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 ρ -approximation algorithm \mathcal{A}' for FTFP.*

Proof. Given an FTFP instance with arbitrary demands, Algorithm \mathcal{A}' works as follows: it solves the LP (3.1) to obtain a fractional optimal solution $(\mathbf{x}^*, \mathbf{y}^*)$, then it constructs instances $\hat{\mathcal{I}}$ and $\dot{\mathcal{I}}$ described above, applies algorithm \mathcal{A} to $\dot{\mathcal{I}}$, 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 $\dot{\mathcal{I}}$ 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 $\rho \cdot \text{cost}(\dot{\mathbf{x}}, \dot{\mathbf{y}})$, because $(\dot{\mathbf{x}}, \dot{\mathbf{y}})$ is feasible for $\dot{\mathcal{I}}$. Thus the cost of \mathcal{A}' is at most

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

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

Two nice consequences are immediate, given demand reduction above.

4.1.1 Reduction from FTFP to FTFL

Given demand reduction, we may assume that we are working with a restricted version of FTFP where every demand r_j is no more than $|\mathbb{F}|$. In this case we can reduce this version of FTFP into FTFL. The reduction simply creates $|\mathbb{F}|$ facilities at each site, and every such facility may be open or close later. Then we have an FTFL instance where every client have a demand r_j and every facility could be open or close. Then any FTFL rounding algorithm can be applied to solve this FTFL instance, and the solution trivially maps into a solution for the corresponding FTFP instance. Moreover, the approximation ratio for FTFL is preserved for FTFP. Given that FTFL has a 1.7245-approximation algorithm by Byrka, Srinivasan and Swamy [7], it is easy to see that FTFP has an approximation algorithm with the same ratio. On the other hand, as we show in Chapter 5, FTFP can be approximated to a ratio of 1.575, a ratio that matches the best LP-based ratio for UFL. So from the standpoint of approximation, FTFP is more amenable than FTFL.

4.1.2 Asymptotic Approximation Ratio for Large Demands

When all demands are large, one would expect the fractional optimal solution to LP (3.1) is very close to an integral solution, and it is reasonable to expect that in this case the fractional solution can be rounded to an integral solution with almost the same cost. We made this intuition a concrete statement, Theorem 10, and provide an affirmative answer to this statement.

Theorem 10 *Given an FTFP instance $(\mathbb{F}, \mathbb{C}, r_j, d_{ij})$ for $i \in \mathbb{F}, j \in \mathbb{C}$, let $m = |\mathbb{F}|$ be the number of sites, and the minimum demand $Q = \min_{j \in \mathbb{C}} r_j$, then there is an approximation algorithm with ratio $1 + O(m/Q)$. In particular, when Q is large compared to m , this ratio approaches 1.*

Proof. We first solve the LP for this instance and obtain optimal fractional solution $(\mathbf{x}^*, \mathbf{y}^*)$, then apply demand reduction to obtain the two instance $\hat{\mathcal{I}}$ and $\dot{\mathcal{I}}$. For the $\hat{\mathcal{I}}$ instance we already have an optimal integral solution, namely $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$. And we now deal with the $\dot{\mathcal{I}}$ instance.

Lemma 8 tells us that \dot{r}_j is no more than $|\mathbb{F}|$ for every client j . So we can solve the $\dot{\mathcal{I}}$ instance by creating m independent UFL instances with the same parameters $\mathbb{F}, \mathbb{C}, d_{ij}$ and all demands are 1. Clearly combining the integral solutions to the m UFL problem would give us an integral solution to the $\dot{\mathcal{I}}$ instance (We might have to remove some redundant connections and facilities, but this only reduces the total cost.). Using an c -approximation algorithm for UFL¹, we can obtain a solution with cost no more than $cm \cdot \text{LP}_{\text{UFL}}$, where

¹We actually need more than that. What we need is that the integral solution for the UFL instance needs to have cost no more than c times the cost of an optimal fractional solution. However, most LP-rounding algorithms do have this property.

LP_{UFL} is the cost of the optimal fractional solution for the UFL instance. On the other hand, it is easy to see that $(\mathbf{x}/Q, \mathbf{y}/Q)$ constitutes a feasible solution to the UFL instance as we have $r_j \geq Q$ for every client j . Therefore, we have an integral solution to instance $\hat{\mathcal{I}}$ with cost at most $(cm/Q)\text{OPT}$ where OPT is the optimal integral solution to the original FTFP instance \mathcal{I} with demands r_j for client j . This solution to $\hat{\mathcal{I}}$, obtained from solutions to m UFL instances, call it S_1 , when combined with the solution to $\hat{\mathcal{I}}$, call it $S_2 = (\hat{\mathbf{x}}, \hat{\mathbf{y}})$, gives a feasible integral solution to the instance \mathcal{I} , and the total cost is no more than

$$\text{cost}(S_1) + \text{cost}(S_2) \leq (cm/Q)\text{OPT} + \text{OPT} = (1 + cm/Q)\text{OPT} = 1 + O(m/Q)\text{OPT}.$$

■

4.2 Adaptive Partition

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 (3.1), 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 ν and μ 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 ν* is $\overline{N}(\nu) = \{\mu \in \overline{\mathbb{F}} : \bar{x}_{\mu\nu} > 0\}$. We will use notation $\nu \in j$ to mean that ν is a demand of client j ; similarly, $\mu \in i$ means that facility μ 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 $C_\nu^{\text{avg}} = \sum_{\mu \in \bar{N}(\nu)} d_{\mu\nu} \bar{x}_{\mu\nu} = \sum_{\mu \in \mathbb{F}} d_{\mu\nu} \bar{x}_{\mu\nu}$. One can think of C_ν^{avg} as the average connection cost of demand ν , if we chose a connection to facility μ with probability $\bar{x}_{\mu\nu}$. In our partitioned fractional solution we guarantee for every ν that $\sum_{\mu \in \mathbb{F}} \bar{x}_{\mu\nu} = 1$.

Some demands in $\bar{\mathbb{C}}$ will be designated as *primary demands* and the set of primary demands will be denoted by P . By definition we have $P \subseteq \bar{\mathbb{C}}$. In addition, we will use the overlap structure between demand neighborhoods to define a mapping that assigns each demand $\nu \in \bar{\mathbb{C}}$ 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. $\sum_{\mu \in \mathbb{F}} \bar{x}_{\mu\nu} = 1$ for each demand $\nu \in \bar{\mathbb{C}}$.
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 ν, ν' of different siblings we have

1. $\overline{N}(\nu) \cap \overline{N}(\nu') = \emptyset$.
2. If ν is assigned to a primary demand κ 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 5.3) is more subtle than that for the 3-approximation

(Section 5.1) and the 1.736-approximation algorithms (Section 5.2), 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 partitioning 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 ν 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 $\bar{x}_{\mu\nu}$ equal to 1, that is condition (PS.1).

For each site i we will initially create one “big” facility μ 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 ν 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 μ . 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 μ after every iteration,

- (c1) For each demand ν 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 μ 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 μ 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}_μ 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 ν for p (Lines 9–10). If $\tilde{N}_1(p)$ overlaps the neighborhood of some existing primary demand κ (if there are multiple such κ 's, pick any of them), we assign ν to κ , and ν acquires all the connection values $\tilde{x}_{\mu p}$ between client p and facility μ 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, ν becomes itself a primary demand and we assign ν 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 ν , that is $\text{conn}(\nu) \stackrel{\text{def}}{=} \sum_{\mu \in \mathbb{F}} \bar{x}_{\mu\nu}$, 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 η 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 ν . Otherwise, $\tilde{x}_{\eta j} > 1 - \text{conn}(\nu)$, in which case we split η so that connecting ν to one of the created copies of η 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 (Pseudocode 1) 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 ν 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 .

Example. We now illustrate our partitioning algorithm with an example, where the FTFP instance has four sites and four clients. The demands are $r_1 = 1$ and $r_2 = r_3 = r_4 = 2$. The facility costs are $f_i = 1$ for all i . The distances are defined as follows: $d_{ii} = 3$ for $i = 1, 2, 3, 4$ and $d_{ij} = 1$ for all $i \neq j$. Solving the LP(3.1), we obtain the fractional solution given in Table 4.1a.

It is easily seen that the fractional solution in Table 4.1a is optimal and complete ($x_{ij}^* > 0$ implies $x_{ij}^* = y_i^*$). The dual optimal solution has all $\alpha_j^* = 4/3$ for $j = 1, 2, 3, 4$.

Now we perform Phase 1, the adaptive partitioning, following the description in Pseudocode 1. To streamline the presentation, we assume that all ties are broken in favor of lower-numbered clients, demands or facilities. First we create one facility at each of the four sites, denoted as $\dot{1}$, $\dot{2}$, $\dot{3}$ and $\dot{4}$ (Line 2–4, Pseudocode 1). We then execute the “while” loop in Line 5 Pseudocode 1. This loop will have seven iterations. Consider the first iteration. In Line 7–8 we compute $\text{tcc}(j)$ for each client $j = 1, 2, 3, 4$ in U . When computing $\tilde{N}_1(2)$, facility $\dot{1}$ will get split into $\dot{1}$ and $\ddot{1}$ with $\bar{y}_{\dot{1}} = 1$ and $\bar{y}_{\ddot{1}} = 1/3$. (This will happen in Line 4–7 of Pseudocode 2.) Then, in Line 9 we will pick client $p = 1$ and create a demand denoted as $1'$ (see Table 4.1b). Since there are no primary demands yet, we make $1'$ a primary demand

x_{ij}^*	1	2	3	4	y_i^*
1	0	$\frac{4}{3}$	$\frac{4}{3}$	$\frac{4}{3}$	$\frac{4}{3}$
2	$\frac{1}{3}$	0	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$
3	$\frac{1}{3}$	$\frac{1}{3}$	0	$\frac{1}{3}$	$\frac{1}{3}$
4	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	0	$\frac{1}{3}$

(a)

$\bar{x}_{\mu\nu}$	1'	2'	2''	3'	3''	4'	4''	\bar{y}_μ
$\dot{1}$	0	1	0	1	0	1	0	1
$\ddot{1}$	0	0	$\frac{1}{3}$	0	$\frac{1}{3}$	0	$\frac{1}{3}$	$\frac{1}{3}$
$\dot{2}$	$\frac{1}{3}$	0	0	0	$\frac{1}{3}$	0	$\frac{1}{3}$	$\frac{1}{3}$
$\dot{3}$	$\frac{1}{3}$	0	$\frac{1}{3}$	0	0	0	$\frac{1}{3}$	$\frac{1}{3}$
$\dot{4}$	$\frac{1}{3}$	0	$\frac{1}{3}$	0	$\frac{1}{3}$	0	0	$\frac{1}{3}$

(b)

Table 4.1: An example of an execution of the partitioning algorithm. (a) An optimal fractional solution x^*, y^* . (b) The partitioned solution. j' and j'' denote the first and second demand of a client j , and \dot{i} and \ddot{i} denote the first and second facility at site i .

with $\bar{N}(1') = \tilde{N}_1(1) = \{\dot{2}, \dot{3}, \dot{4}\}$. Notice that client 1 is exhausted after this iteration and U becomes $\{2, 3, 4\}$.

In the second iteration we compute $\text{tcc}(j)$ for $j = 2, 3, 4$ and pick client $p = 2$, from which we create a new demand $2'$. We have $\tilde{N}_1(2) = \{\dot{1}\}$, which is disjoint from $\bar{N}(1')$. So we create a demand $2'$ and make it primary, and set $\bar{N}(2') = \{\dot{1}\}$. In the third iteration we compute $\text{tcc}(j)$ for $j = 2, 3, 4$ and again we pick client $p = 2$. Since $\tilde{N}_1(2) = \{\ddot{1}, \dot{3}, \dot{4}\}$ overlaps with $\bar{N}(1')$, we create a demand $2''$ and assign it to $1'$. We also set $\bar{N}(2'') = \bar{N}(1') \cap \tilde{N}(2) = \{\dot{3}, \dot{4}\}$. After this iteration client 2 is exhausted and we have $U = \{3, 4\}$.

In the fourth iteration we compute $\text{tcc}(j)$ for client $j = 3, 4$. We pick $p = 3$ and

create demand $3'$. Since $\tilde{N}_1(3) = \{\dot{1}\}$ overlaps $\bar{N}(2')$, we assign $3'$ to $2'$ and set $\bar{N}(3') = \{\dot{1}\}$. In the fifth iteration we compute $\text{tcc}(j)$ for client $j = 3, 4$ and pick $p = 3$ again. At this time $\tilde{N}_1(3) = \{\ddot{1}, \dot{2}, \dot{4}\}$, which overlaps with $\bar{N}(1')$. So we create a demand $3''$ and assign it to $1'$, as well as set $\bar{N}(3'') = \{\dot{2}, \dot{4}\}$.

In the last two iterations we will pick client $p = 4$ twice and create demands $4'$ and $4''$. For $4'$ we have $\tilde{N}_1(4) = \{\dot{1}\}$ so we assign $4'$ to $2'$ and set $\bar{N}(4') = \{\dot{1}\}$. For $4''$ we have $\tilde{N}_1(4) = \{\ddot{1}, \dot{2}, \dot{3}\}$ and we assign it to $1'$, as well as set $\bar{N}(4'') = \{\dot{2}, \dot{3}\}$.

Now that all clients are exhausted we perform Phase 2, the augmenting phase, to construct a fractional solution in which all demands have total connection value equal to 1. We iterate through each of the seven demands created, that is $1', 2', 2'', 3', 3'', 4', 4''$. $1'$ and $2'$ already have neighborhoods with total connection value of 1, so nothing will change in the first two iterations. $2''$ has $\dot{3}, \dot{4}$ in its neighborhood, with total connection value of $2/3$, and $\tilde{N}(2) = \{\ddot{1}\}$ at this time, so we add $\ddot{1}$ into $\bar{N}(2'')$ to make $\bar{N}(2'') = \{\ddot{1}, \dot{3}, \dot{4}\}$ and now $2''$ has total connection value of 1. Similarly, $3''$ and $4''$ each get $\ddot{1}$ added to their neighborhood and end up with total connection value of 1. The other two demands, namely $3'$ and $4'$, each have $\dot{1}$ in its neighborhood so each of them has already its total connection value equal 1. This completes Phase 2.

The final partitioned fractional solution is given in Table 4.1b. We have created a total of five facilities $\dot{1}, \ddot{1}, \dot{2}, \dot{3}, \dot{4}$, and seven demands, $1', 2', 2'', 3', 3'', 4', 4''$. It can be verified that all the stated properties are satisfied.

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 μ into μ' and μ'' , the sum of $\bar{y}_{\mu'}$ and $\bar{y}_{\mu''}$ is equal to the old value of \bar{y}_{μ} .

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 μ 's for which such $\kappa \in P$ exists, and denote this κ by κ_{μ} . Then, using conditions (CO) and (PS.3), we have $\sum_{\mu \in i} \sum_{\kappa \in P} \bar{x}_{\mu\kappa} = \sum_{\mu \in K} \bar{x}_{\mu\kappa_{\mu}} = \sum_{\mu \in K} \bar{y}_{\mu} \leq \sum_{\mu \in i} \bar{y}_{\mu} = y_i^*$.

Property (PD.3(a)) follows from the way the algorithm assigns primary demands. When demand ν of client p is assigned to a primary demand κ 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 13.

Finally we argue that the properties in group (SI) hold. (SI.1) is easy, since for any client j , each facility μ is added to the neighborhood of at most one demand $\nu \in j$, by setting $\bar{x}_{\mu\nu}$ to \bar{y}_μ , while other siblings ν' of ν 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 11.

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

Proof. Let ν_1, \dots, ν_{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 κ_q the primary demand that ν_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 $\bar{N}(\nu_s) \cap \bar{N}(\kappa_q) = \emptyset$ for any $s \neq q$, that is Property (SI.2) holds right after Phase 1.

After Phase 1 all neighborhoods $\bar{N}(\kappa_s)$, $s = 1, \dots, r_j$ have already been fixed and they do not change in Phase 2. None of the facilities in $\tilde{N}(j)$ appear in any of $\bar{N}(\kappa_s)$ 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 $\bar{N}(\nu)$, for some $\nu \in j$ (Line 19–21 of Pseudocode 1), all the required disjointness conditions will be preserved. ■

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

Lemma 12 *Let η and ν be two demands, with η created no later than ν , and let $j \in \mathbb{C}$ be a client that is not exhausted when ν 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 η 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 μ_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_ν^{avg} if $\nu \in j$. Indeed, each of these quantities involves some later values $\tilde{x}_{\mu j}$, where μ could be one of the facilities μ_s or a new facility obtained from splitting. For each s , however, the sum of all values $\tilde{x}_{\mu j}$,

over the facilities μ that were split from μ_s , cannot exceed the value $\tilde{x}_{\mu_s j}$ at the time when η 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_ν^{avg} (for $\nu \in j$) correspond to some choice of the z_s variables (adding up to 1), and the lemma follows. ■

Lemma 13 *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 κ is a primary demand of client p , then C_κ^{avg} is equal to $\text{tcc}(p)$ computed when κ 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 12. ■

We have thus proved that all properties (PS), (CO), (PD) and (SI) hold for our partitioned fractional solution $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$. 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 5.1) and the 1.736-approximation algorithm (Section 5.2), the first phase of the algorithm is exactly the same partition process as described above. However, the 1.575-approximation algorithm (Section 5.3) 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.

Pseudocode 1 Algorithm: Adaptive Partitioning

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

```

1:  $\tilde{\mathbf{r}} \leftarrow \mathbf{r}, U \leftarrow \mathbb{C}, \overline{\mathbb{F}} \leftarrow \emptyset, \overline{\mathbb{C}} \leftarrow \emptyset, P \leftarrow \emptyset$  ▷ Phase 1

2: for each site  $i \in \mathbb{F}$  do

3:   create a facility  $\mu$  at  $i$  and add  $\mu$  to  $\overline{\mathbb{F}}$ 

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:      $\tilde{N}_1(j) \leftarrow \text{NEARESTUNITCHUNK}(j, \overline{\mathbb{F}}, \tilde{\mathbf{x}}, \bar{\mathbf{x}}, \bar{\mathbf{y}})$  ▷ 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  $\tilde{N}_1(p) \cap \overline{N}(\kappa) \neq \emptyset$  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  $\mu \in \tilde{N}_1(p) \cap \overline{N}(\kappa)$ 

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:       $\overline{\mathbb{C}} \leftarrow \overline{\mathbb{C}} \cup \{\nu\}, \tilde{r}_p \leftarrow \tilde{r}_p - 1$ 

18:      if  $\tilde{r}_p = 0$  then  $U \leftarrow U \setminus \{p\}$ 

19: for each client  $j \in \mathbb{C}$  do ▷ Phase 2

20:   for each demand  $\nu \in j$  do ▷ each client  $j$  has  $r_j$  demands

21:     if  $\sum_{\mu \in \overline{N}(\nu)} \bar{x}_{\mu\nu} < 1$  then  $\text{AUGMENTTOUNIT}(\nu, j, \overline{\mathbb{F}}, \tilde{\mathbf{x}}, \bar{\mathbf{x}}, \bar{\mathbf{y}})$  ▷ see Pseudocode 2

```

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  $\nu \in \bar{\mathbb{C}}$  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  $j' \in \bar{\mathbb{C}}$  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)

```

Chapter 5

LP-rounding Algorithms

5.1 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 [15].

Algorithm EGUP. At a high level, we would open exactly one facility for each primary demand κ , 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}_μ)

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 κ .

In our description above, the algorithm is presented as a randomized algorithm. It can be de-randomized 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 [9], 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 14 *The expectation of facility cost F_{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

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

where the inequality follows from Property (PD.2). ■

Lemma 15 *The expectation of connection cost C_{EGUP} of our solution is at most $C^* + 2 \cdot \text{LP}^*$.*

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

Consider a non-primary demand ν assigned to a primary demand $\kappa \in P$. Let μ be any facility in $\bar{N}(\nu) \cap \bar{N}(\kappa)$. Since μ 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 ν to its facility $\phi(\kappa)$ is

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

where the second inequality follows from Property (PD.3(b)). From the definition of C_{ν}^{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

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

completing the proof of the lemma. ■

Theorem 16 *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 14 and Lemma 15 imply that the total cost is at most $F^* + C^* + 2 \cdot \text{LP}^* = 3 \cdot \text{LP}^* \leq 3 \cdot \text{OPT}$. ■

5.2 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 could be far from the optimum, since we connect a non-primary demand to a facility in the neighborhood of its assigned primary demand and then estimate the distance using the triangle inequality. The basic idea to improve the estimate of the connection cost, following the approach of Chudak and Shmoys [9], is to connect each non-primary demand to its nearest neighbor when one is available, and to only use the facility opened by its assigned primary demand when none of its neighbors is open.

Algorithm ECHS. As before, the algorithm starts by solving the linear program and applying the adaptive partitioning algorithm described in Section 4.2 to obtain a partitioned solution (\bar{x}, \bar{y}) . Then we apply the rounding process to compute an integral solution (see Pseudocode 3).

We start, as before, by opening exactly one facility $\phi(\kappa)$ in the neighborhood of each primary demand κ (Line 2). For any non-primary demand ν assigned to κ , we refer to $\phi(\kappa)$ as the *target* facility of ν . In Algorithm EGUP, ν was connected to $\phi(\kappa)$, but in Algorithm ECHS we may be able to find an open facility in ν 's neighborhood and connect ν to this facility. Specifically, the two changes in the algorithm are as follows:

- (1) Each facility μ that is not in the neighborhood of any primary demand is opened, independently, with probability \bar{y}_μ (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 ν . 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 κ is connected to the only facility $\phi(\kappa)$ opened in its neighborhood, as before (Line 3). For a non-primary demand ν , if its neighborhood $\overline{N}(\nu)$ has an open facility, we connect ν to the closest open facility in $\overline{N}(\nu)$ (Line 8). Otherwise, we connect ν to its target facility (Line 10).

Pseudocode 3 Algorithm ECHS: Constructing Integral Solution

```

1: for each  $\kappa \in P$  do
2:   choose one  $\phi(\kappa) \in \overline{N}(\kappa)$ , with each  $\mu \in \overline{N}(\kappa)$  chosen as  $\phi(\kappa)$  with probability  $\bar{y}_\mu$ 
3:   open  $\phi(\kappa)$  and connect  $\kappa$  to  $\phi(\kappa)$ 
4: for each  $\mu \in \overline{\mathbb{F}} - \bigcup_{\kappa \in P} \overline{N}(\kappa)$  do
5:   open  $\mu$  with probability  $\bar{y}_\mu$  (independently)
6: for each non-primary demand  $\nu \in \overline{\mathbb{C}}$  do
7:   if any facility in  $\overline{N}(\nu)$  is open then
8:     connect  $\nu$  to the nearest open facility in  $\overline{N}(\nu)$ 
9:   else
10:    connect  $\nu$  to  $\phi(\kappa)$  where  $\kappa$  is  $\nu$ 's assigned 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 ν and ν' be two different sibling demands of client j and let their assigned primary demands be κ and κ' respectively. Due to (SI.2) we know $\kappa \neq \kappa'$. From (SI.1) we have $\overline{N}(\nu) \cap \overline{N}(\nu') = \emptyset$. From (SI.2), we have $\overline{N}(\nu) \cap \overline{N}(\kappa') = \emptyset$ and $\overline{N}(\nu') \cap \overline{N}(\kappa) = \emptyset$. From (PD.1) we have $\overline{N}(\kappa) \cap \overline{N}(\kappa') = \emptyset$. It

follows that $(\overline{N}(\nu) \cup \overline{N}(\kappa)) \cap (\overline{N}(\nu') \cup \overline{N}(\kappa')) = \emptyset$. Since the algorithm connects ν to some facility in $\overline{N}(\nu) \cup \overline{N}(\kappa)$ and ν' to some facility in $\overline{N}(\nu') \cup \overline{N}(\kappa')$, ν and ν' 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

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

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

To bound the connection cost, we adapt an argument of Chudak and Shmoys [9]. Consider a demand ν and denote by C_{ν} the random variable representing the connection cost for ν . Our goal now is to estimate $\mathbb{E}[C_{\nu}]$, the expected value of C_{ν} . Demand ν can either get connected directly to some facility in $\overline{N}(\nu)$ or indirectly to its target facility $\phi(\kappa) \in \overline{N}(\kappa)$, where κ is the primary demand to which ν is assigned. We will analyze these two cases separately.

In our analysis, in this section and the next one, we will use notation

$$D(A, \sigma) = \sum_{\mu \in A} d_{\mu\sigma} \bar{y}_{\mu} / \sum_{\mu \in A} \bar{y}_{\mu}$$

for the average distance between a demand σ and a set A of facilities. Note that, in particular, we have $C_{\nu}^{\text{avg}} = D(\overline{N}(\nu), \nu)$.

We first estimate the expected cost $d_{\phi(\kappa)\nu}$ of the indirect connection. Let Λ^ν denote the event that some facility in $\overline{N}(\nu)$ is opened. Then

$$\mathbb{E}[C_\nu \mid \neg\Lambda^\nu] = \mathbb{E}[d_{\phi(\kappa)\nu} \mid \neg\Lambda^\nu] = D(\overline{N}(\kappa) \setminus \overline{N}(\nu), \nu). \quad (5.1)$$

Note that $\neg\Lambda^\nu$ implies that $\overline{N}(\kappa) \setminus \overline{N}(\nu) \neq \emptyset$, since $\overline{N}(\kappa)$ contains exactly one open facility, namely $\phi(\kappa)$.

Lemma 17 *Let ν be a demand assigned to a primary demand κ , and assume that $\overline{N}(\kappa) \setminus \overline{N}(\nu) \neq \emptyset$. Then*

$$\mathbb{E}[C_\nu \mid \neg\Lambda^\nu] \leq C_\nu^{\text{avg}} + 2\alpha_\nu^*.$$

Proof. By (5.1), we need to show that $D(\overline{N}(\kappa) \setminus \overline{N}(\nu), \nu) \leq C_\nu^{\text{avg}} + 2\alpha_\nu^*$. There are two cases to consider.

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

$\mu \in \overline{N}(\kappa) \setminus \overline{N}(\nu)$, 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 $\mu \in \overline{N}(\kappa) \setminus \overline{N}(\nu)$, it follows that $D(\overline{N}(\kappa) \setminus \overline{N}(\nu), \nu) \leq C_\nu^{\text{avg}} + 2\alpha_\nu^*$.

Case 2: Every $\mu' \in \overline{N}(\kappa) \cap \overline{N}(\nu)$ has $d_{\mu'\kappa} > C_\kappa^{\text{avg}}$. Since $C_\kappa^{\text{avg}} = D(\overline{N}(\kappa), \kappa)$, this implies

that $D(\overline{N}(\kappa) \setminus \overline{N}(\nu), \kappa) \leq C_\kappa^{\text{avg}}$. Therefore, choosing an arbitrary $\mu' \in \overline{N}(\kappa) \cap \overline{N}(\nu)$,

we obtain

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

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. ■

We now continue our estimation of the connection cost. The next step of our analysis is to show that

$$\mathbb{E}[C_\nu] \leq C_\nu^{\text{avg}} + \frac{2}{e}\alpha_\nu^*. \quad (5.2)$$

The argument is divided into three cases. The first, easy case is when ν is a primary demand κ . According to the algorithm (see Pseudocode 3, Line 2), we have $C_\kappa = d_{\mu\kappa}$ with probability \bar{y}_μ , for $\mu \in \bar{N}(\kappa)$. Therefore $\mathbb{E}[C_\kappa] = C_\kappa^{\text{avg}}$, so (5.2) holds.

Next, we consider a non-primary demand ν . Let κ be the primary demand that ν is assigned to. We first deal with the sub-case when $\bar{N}(\kappa) \setminus \bar{N}(\nu) = \emptyset$, which is the same as $\bar{N}(\kappa) \subseteq \bar{N}(\nu)$. Property (CO) implies that $\bar{x}_{\mu\nu} = \bar{y}_\mu = \bar{x}_{\mu\kappa}$ for every $\mu \in \bar{N}(\kappa)$, so we have $\sum_{\mu \in \bar{N}(\kappa)} \bar{x}_{\mu\nu} = \sum_{\mu \in \bar{N}(\kappa)} \bar{x}_{\mu\kappa} = 1$, due to (PS.1). On the other hand, we have $\sum_{\mu \in \bar{N}(\nu)} \bar{x}_{\mu\nu} = 1$, and $\bar{x}_{\mu\nu} > 0$ for all $\mu \in \bar{N}(\nu)$. Therefore $\bar{N}(\kappa) = \bar{N}(\nu)$ and C_ν has exactly the same distribution as C_κ . So this case reduces to the first case, namely we have $\mathbb{E}[C_\nu] = C_\nu^{\text{avg}}$, and (5.2) holds.

The last, and only non-trivial case is when $\bar{N}(\kappa) \setminus \bar{N}(\nu) \neq \emptyset$. We handle this case in the following lemma.

Lemma 18 *Assume that $\bar{N}(\kappa) \setminus \bar{N}(\nu) \neq \emptyset$. Then the expected connection cost of ν , conditioned on the event that at least one of its neighbor opens, satisfies*

$$\mathbb{E}[C_\nu \mid \Lambda^\nu] \leq C_\nu^{\text{avg}}.$$

Proof. The proof is similar to an analogous result in [9, 5]. For the sake of completeness we sketch here a simplified argument, adapted to our terminology and nota-

tion. The idea is to consider a different random process that is easier to analyze and whose expected connection cost is not better than that in the algorithm.

We partition $\overline{N}(\nu)$ into groups G_1, \dots, G_k , where two different facilities μ and μ' are put in the same G_s , where $s \in \{1, \dots, k\}$, if they both belong to the same set $\overline{N}(\kappa)$ for some primary demand κ . If some μ is not a neighbor of any primary demand, then it constitutes a singleton group. For each s , let $\bar{d}_s = D(G_s, \nu)$ be the average distance from ν to G_s . Assume that G_1, \dots, G_k are ordered by nondecreasing average distance to ν , that is $\bar{d}_1 \leq \bar{d}_2 \leq \dots \leq \bar{d}_k$. For each group G_s , we select it, independently, with probability $g_s = \sum_{\mu \in G_s} \bar{y}_\mu$. For each selected group G_s , we open exactly one facility in G_s , where each $\mu \in G_s$ is opened with probability $\bar{y}_\mu / \sum_{\eta \in G_s} \bar{y}_\eta$.

So far, this process is the same as that in the algorithm (if restricted to $\overline{N}(\nu)$). However, we connect ν in a slightly different way, by choosing the smallest s for which G_s was selected and connecting ν to the open facility in G_s . This can only increase our expected connection cost, assuming that at least one facility in $\overline{N}(\nu)$ opens, so

$$\begin{aligned} \mathbb{E}[C_\nu \mid \Lambda^\nu] &\leq \frac{1}{\mathbb{P}[\Lambda^\nu]} (\bar{d}_1 g_1 + \bar{d}_2 g_2 (1 - g_1) + \dots + \bar{d}_k g_k (1 - g_1)(1 - g_2) \dots (1 - g_{k-1})) \\ &\leq \frac{1}{\mathbb{P}[\Lambda^\nu]} \cdot \sum_{s=1}^k \bar{d}_s g_s \cdot \left(\sum_{t=1}^k g_t \prod_{z=1}^{t-1} (1 - g_z) \right) \end{aligned} \quad (5.3)$$

$$= \sum_{s=1}^k \bar{d}_s g_s \quad (5.4)$$

$$= C_\nu^{\text{avg}}. \quad (5.5)$$

The proof for inequality (5.3) is given in A.2 (note that $\sum_{s=1}^k g_s = 1$), equality (5.4) follows from $\mathbb{P}[\Lambda^\nu] = 1 - \prod_{t=1}^k (1 - g_t) = \sum_{t=1}^k g_t \prod_{z=1}^{t-1} (1 - g_z)$, and (5.5) follows from the definition of the distances \bar{d}_s , probabilities g_s , and simple algebra. ■

Next, we show an estimate on the probability that none of ν 's neighbors is opened by the algorithm.

Lemma 19 *The probability that none of ν 's neighbors is opened satisfies $\mathbb{P}[\neg\Lambda^\nu] \leq 1/e$.*

Proof. We use the same partition of $\overline{N}(\nu)$ into groups G_1, \dots, G_k as in the proof of Lemma 18. Denoting by g_s the probability that a group G_s is selected (and thus that it has an open facility), we have

$$\mathbb{P}[\neg\Lambda^\nu] = \prod_{s=1}^k (1 - g_s) \leq e^{-\sum_{s=1}^k g_s} = e^{-\sum_{\mu \in \overline{N}(\nu)} \bar{y}_\mu} = \frac{1}{e}.$$

In this derivation, we first use that $1 - x \leq e^{-x}$ holds for all x , the second equality follows from $\sum_{s=1}^k g_s = \sum_{\mu \in \overline{N}(\nu)} \bar{y}_\mu$ and the last equality follows from $\sum_{\mu \in \overline{N}(\nu)} \bar{y}_\mu = 1$. ■

We are now ready to estimate the unconditional expected connection cost of ν (in the case when $\overline{N}(\kappa) \setminus \overline{N}(\nu) \neq \emptyset$) as follows,

$$\begin{aligned} \mathbb{E}[C_\nu] &= \mathbb{E}[C_\nu \mid \Lambda^\nu] \cdot \mathbb{P}[\Lambda^\nu] + \mathbb{E}[C_\nu \mid \neg\Lambda^\nu] \cdot \mathbb{P}[\neg\Lambda^\nu] \\ &\leq C_\nu^{\text{avg}} \cdot \mathbb{P}[\Lambda^\nu] + (C_\nu^{\text{avg}} + 2\alpha_\nu^*) \cdot \mathbb{P}[\neg\Lambda^\nu] \end{aligned} \tag{5.6}$$

$$\begin{aligned} &= C_\nu^{\text{avg}} + 2\alpha_\nu^* \cdot \mathbb{P}[\neg\Lambda^\nu] \\ &\leq C_\nu^{\text{avg}} + \frac{2}{e} \cdot \alpha_\nu^*. \end{aligned} \tag{5.7}$$

In the above derivation, inequality (5.6) follows from Lemmas 17 and 18, and inequality (5.7) follows from Lemma 19.

We have thus shown that the bound (5.2) holds in all three cases. Summing over all demands ν of a client j , we can now bound the expected connection cost of client j :

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

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

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

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

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

Summarizing, we obtain the main result of this section.

Theorem 20 *Algorithm ECHS is a $(1 + 2/e)$ -approximation algorithm for FTFP.*

5.3 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.* [6].

Recall that in Section 5.2 we showed how to compute an integral solution with facility cost bounded by F^* and connection cost bounded by $C^* + 2/e \cdot \text{LP}^*$. Thus, while our facility cost does not exceed the optimal fractional facility cost, our connection cost is significantly larger than the connection cost in the optimal fractional solution. A natural idea is to balance these two ratios by reducing the connection cost at the expense of the facility cost. One way to do this would be to increase the probability of opening facilities, from \bar{y}_μ (used in Algorithm ECHS) to, say, $\gamma \bar{y}_\mu$, for some $\gamma > 1$. This increases the expected

facility cost by a factor of γ but, as it turns out, it also reduces the probability that an indirect connection occurs for a non-primary demand to $1/e^\gamma$ (from the previous value $1/e$ in ECHS). As a consequence, for each primary demand κ , the new algorithm will select a facility to open from the nearest facilities μ in $\overline{N}(\kappa)$ such that the connection values $\bar{x}_{\mu\nu}$ sum up to $1/\gamma$, instead of 1 as in Algorithm ECHS. It is easily seen that this will improve the estimate on connection cost for primary demands. These two changes, along with a more refined analysis, are the essence of the approach in [6], expressed in our terminology.

Our approach can be thought of as a combination of the above ideas 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 disjoint parts, and with restrictions on how parts from different demands can overlap.

We begin by describing properties that our partitioned fractional solution (\bar{x}, \bar{y}) needs to satisfy. Assume that γ is some constant such that $1 < \gamma < 2$. As mentioned earlier, the neighborhood $\overline{N}(\nu)$ of each demand ν will be divided into two disjoint parts. The first part, called the *close neighborhood* and denoted $\overline{N}_{\text{cls}}(\nu)$, contains the facilities in $\overline{N}(\nu)$ nearest to ν with the total connection value equal $1/\gamma$, that is $\sum_{\mu \in \overline{N}_{\text{cls}}(\nu)} \bar{x}_{\mu\nu} = 1/\gamma$. The second part, called the *far neighborhood* and denoted $\overline{N}_{\text{far}}(\nu)$, contains the remaining facilities in $\overline{N}(\nu)$ (so $\sum_{\mu \in \overline{N}_{\text{far}}(\nu)} \bar{x}_{\mu\nu} = 1 - 1/\gamma$). We restate these definitions formally below in Property (NB). Recall that for any set A of facilities and a demand ν , by $D(A, \nu)$ we denote the average distance between ν and the facilities in A , that is $D(A, \nu) = \sum_{\mu \in A} d_{\mu\nu} \bar{y}_\mu / \sum_{\mu \in A} \bar{y}_\mu$.

We will use notations $C_{\text{cls}}^{\text{avg}}(\nu) = D(\overline{N}_{\text{cls}}(\nu), \nu)$ and $C_{\text{far}}^{\text{avg}}(\nu) = D(\overline{N}_{\text{far}}(\nu), \nu)$ for the average distances from ν to its close and far neighborhoods, respectively. By the definition of these sets and the completeness property (CO), these distances can be expressed as

$$C_{\text{cls}}^{\text{avg}}(\nu) = \gamma \sum_{\mu \in \overline{N}_{\text{cls}}(\nu)} d_{\mu\nu} \bar{x}_{\mu\nu} \quad \text{and} \quad C_{\text{far}}^{\text{avg}}(\nu) = \frac{\gamma}{\gamma - 1} \sum_{\mu \in \overline{N}_{\text{far}}(\nu)} d_{\mu\nu} \bar{x}_{\mu\nu}.$$

We will also use notation $C_{\text{cls}}^{\text{max}}(\nu) = \max_{\mu \in \overline{N}_{\text{cls}}(\nu)} d_{\mu\nu}$ for the maximum distance from ν to its close neighborhood. The average distance from a demand ν to its overall neighborhood $\overline{N}(\nu)$ is denoted as $C^{\text{avg}}(\nu) = D(\overline{N}(\nu), \nu) = \sum_{\mu \in \overline{N}(\nu)} d_{\mu\nu} \bar{x}_{\mu\nu}$. It is easy to see that

$$C^{\text{avg}}(\nu) = \frac{1}{\gamma} C_{\text{cls}}^{\text{avg}}(\nu) + \frac{\gamma - 1}{\gamma} C_{\text{far}}^{\text{avg}}(\nu). \quad (5.8)$$

Our partitioned solution $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$ must satisfy the same partitioning and completeness properties as before, namely properties (PS) and (CO) in Section 4.2. In addition, it must satisfy a new neighborhood property (NB) and modified properties (PD') and (SI'), listed below.

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

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

Note that the first two conditions, together with (PS.1), imply that $\sum_{\mu \in \overline{N}_{\text{far}}(\nu)} \bar{x}_{\mu\nu} = 1 - 1/\gamma$. When defining $\overline{N}_{\text{cls}}(\nu)$, in case of ties, which can occur when some facilities

in $\overline{N}(\nu)$ are at the same distance from ν , we use a tie-breaking rule that is explained in the proof of Lemma 21 (the only place where the rule is needed).

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

1. For any two different primary demands $\kappa, \kappa' \in P$ we have $\overline{N}_{\text{cls}}(\kappa) \cap \overline{N}_{\text{cls}}(\kappa') = \emptyset$.
2. For each site $i \in \mathbb{F}$, $\sum_{\kappa \in P} \sum_{\mu \in i \cap \overline{N}_{\text{cls}}(\kappa)} \bar{x}_{\mu\kappa} \leq y_i^*$. In the summation, as before, we overload notation i to stand for the set of facilities created on site i .
3. Each demand $\nu \in \overline{\mathbb{C}}$ is assigned to one primary demand $\kappa \in P$ such that
 - (a) $\overline{N}_{\text{cls}}(\nu) \cap \overline{N}_{\text{cls}}(\kappa) \neq \emptyset$, 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 $\nu, \nu' \in \overline{\mathbb{C}}$ of different siblings we have

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

Modified adaptive partitioning. To obtain a fractional solution with the above properties, we employ a modified adaptive partitioning algorithm. As in Section 4.2, 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 $\bar{x}_{\mu\nu}$ so that the total connection value of each demand ν is 1. As the partitioning algorithm proceeds, for any demand ν , $\overline{N}(\nu)$ denotes the set of facilities with $\bar{x}_{\mu\nu} > 0$; hence the notation $\overline{N}(\nu)$ actually represents

a dynamic set which gets fixed once the partitioning algorithm concludes both Phase 2. On the other hand, $\overline{N}_{\text{cls}}(\nu)$ and $\overline{N}_{\text{far}}(\nu)$ refer to the close and far neighborhoods at the time when $\overline{N}(\nu)$ is fixed.

Similar to the algorithm in Section 4.2, Phase 1 runs in iterations. Fix some iteration and 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 μ such that $\tilde{x}_{\mu j} > 0$. Order the facilities in this set as $\tilde{N}(j) = \{\mu_1, \dots, \mu_q\}$ with non-decreasing distance from j , that is $d_{\mu_1 j} \leq d_{\mu_2 j} \leq \dots \leq d_{\mu_q j}$. Without loss of generality, there is an index l for which $\sum_{s=1}^l \tilde{x}_{\mu_s j} = 1/\gamma$, since we can always split one facility to achieve this. Then we define $\tilde{N}_{\text{cls}}(j) = \{\mu_1, \dots, \mu_l\}$. (Unlike close neighborhoods of demands, $\tilde{N}_{\text{cls}}(j)$ can vary over time.)

We also use notation

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

When the iteration starts, we first find a not-yet-exhausted client p that minimizes the value of $\text{tcc}_{\text{cls}}(p) + \text{dmax}_{\text{cls}}(p)$ and create a new demand ν for p . Now we have two cases:

Case 1: $\tilde{N}_{\text{cls}}(p) \cap \overline{N}(\kappa) \neq \emptyset$ for some existing primary demand $\kappa \in P$. In this case we assign ν to κ . As before, if there are multiple such κ , we pick any of them. We also fix $\bar{x}_{\mu\nu} \leftarrow \tilde{x}_{\mu p}$ and $\tilde{x}_{\mu p} \leftarrow 0$ for each $\mu \in \tilde{N}(p) \cap \overline{N}(\kappa)$. Note that although we check for overlap between $\tilde{N}_{\text{cls}}(p)$ and $\overline{N}(\kappa)$, the facilities we actually move into $\overline{N}(\nu)$ include all facilities in the intersection of $\tilde{N}(p)$, a bigger set, with $\overline{N}(\kappa)$.

At this time, the total connection value between ν and $\mu \in \overline{N}(\nu)$ is at most $1/\gamma$, since

$\sum_{\mu \in \overline{N}(\kappa)} \bar{y}_{\mu} = 1/\gamma$ (this follows from the definition of neighborhoods for new primary

demands in Case 2 below) and we have $\overline{N}(\nu) \subseteq \overline{N}(\kappa)$ at this point. Later in Phase 2 we will add additional facilities from $\tilde{N}(p)$ to $\overline{N}(\nu)$ to make ν 's total connection value equal to 1.

Case 2: $\tilde{N}_{\text{cls}}(p) \cap \overline{N}(\kappa) = \emptyset$ for all existing primary demands $\kappa \in P$. In this case we make

ν a primary demand (that is, add it to P) and assign it to itself. We then move the facilities from $\tilde{N}_{\text{cls}}(p)$ to $\overline{N}(\nu)$, that is for $\mu \in \tilde{N}_{\text{cls}}(p)$ we set $\bar{x}_{\mu\nu} \leftarrow \tilde{x}_{\mu p}$ and $\tilde{x}_{\mu p} \leftarrow 0$.

It is easy to see that the total connection value of ν to $\overline{N}(\nu)$ is now exactly $1/\gamma$, that is $\sum_{\mu \in \overline{N}(\nu)} \bar{y}_\mu = 1/\gamma$. Moreover, facilities remaining in $\tilde{N}(p)$ are all farther away from ν than those in $\overline{N}(\nu)$. As we add only facilities from $\tilde{N}(p)$ to $\overline{N}(\nu)$ in Phase 2, the final $\overline{N}_{\text{cls}}(\nu)$ contains the same set of facilities as the current set $\overline{N}(\nu)$. (More precisely, $\overline{N}_{\text{cls}}(\nu)$ consists of the facilities that either are currently in $\overline{N}(\nu)$ or were obtained from splitting the facilities currently in $\overline{N}(\nu)$.)

Once all clients are exhausted, that is, each client j has r_j demands created, Phase 1 concludes. We then run Phase 2, the augmenting phase, following the same steps as in Section 4.2. For each client j and each demand $\nu \in j$ with total connection value to $\overline{N}(\nu)$ less than 1 (that is, $\sum_{\mu \in \overline{N}(\nu)} \bar{x}_{\mu\nu} < 1$), we use our `AUGMENTTOUNIT()` procedure to add additional facilities (possibly split, if necessary) from $\tilde{N}(j)$ to $\overline{N}(\nu)$ to make the total connection value between ν and $\overline{N}(\nu)$ equal 1.

This completes the description of the partitioning algorithm. Summarizing, for each client $j \in \mathbb{C}$ we created r_j demands on the same point as j , and we created a number of facilities at each site $i \in \mathbb{F}$. Thus computed sets of demands and facilities are denoted

$\overline{\mathbb{C}}$ and $\overline{\mathbb{F}}$, respectively. For each facility $\mu \in i$ we defined its fractional opening value \bar{y}_μ , $0 \leq \bar{y}_\mu \leq 1$, and for each demand $\nu \in j$ we defined its fractional connection value $\bar{x}_{\mu\nu} \in \{0, \bar{y}_\mu\}$. The connections with $\bar{x}_{\mu\nu} > 0$ define the neighborhood $\overline{N}(\nu)$. The facilities in $\overline{N}(\nu)$ that are closest to ν and have total connection value from ν equal $1/\gamma$ form the close neighborhood $\overline{N}_{\text{cls}}(\nu)$, while the remaining facilities in $\overline{N}(\nu)$ form the far neighborhood $\overline{N}_{\text{far}}(\nu)$. It remains to show that this partitioning satisfies all the desired properties.

Correctness of partitioning. We now argue that our partitioned fractional solution (\bar{x}, \bar{y}) satisfies all the stated properties. Properties (PS), (CO) and (NB) are directly enforced by the algorithm.

(PD'.1) holds because for each primary demand $\kappa \in p$, $\overline{N}_{\text{cls}}(\kappa)$ is the same set as $\tilde{N}_{\text{cls}}(p)$ at the time when κ was created, and $\tilde{N}_{\text{cls}}(p)$ is removed from $\tilde{N}(p)$ right after this step. Further, the partitioning algorithm makes κ a primary demand only if $\tilde{N}_{\text{cls}}(p)$ is disjoint from the set $\overline{N}(\kappa')$ of all existing primary demands κ' at that iteration, but these neighborhoods are the same as the final close neighborhoods $\overline{N}_{\text{cls}}(\kappa')$.

The justification of (PD'.2) is similar to that for (PD.2) from Section 4.2. All close neighborhoods of primary demands are disjoint, due to (PD'.1), so each facility $\mu \in i$ can appear in at most one $\overline{N}_{\text{cls}}(\kappa)$, for some $\kappa \in P$. Condition (CO) implies that $\bar{y}_\mu = \bar{x}_{\mu\kappa}$ for $\mu \in \overline{N}_{\text{cls}}(\kappa)$. As a result, the summation on the left-hand side is not larger than $\sum_{\mu \in i} \bar{y}_\mu = y_i^*$.

Regarding (PD'.3(a)), at first glance this property seems to follow directly from the algorithm, as we only assign a demand ν to a primary demand κ when $\overline{N}(\nu)$ at that iteration overlaps with $\overline{N}(\kappa)$ (which is equal to the final value of $\overline{N}_{\text{cls}}(\kappa)$). However, it is

a little more subtle, as the final $\overline{N}_{\text{cls}}(\nu)$ may contain facilities added to $\overline{N}(\nu)$ in Phase 2. Those facilities may turn out to be closer to ν than some facilities in $\overline{N}(\kappa) \cap \tilde{N}(j)$ (not $\tilde{N}_{\text{cls}}(j)$) that we added to $\overline{N}(\nu)$ in Phase 1. If the final $\overline{N}_{\text{cls}}(\nu)$ consists only of facilities added in Phase 2, we no longer have the desired overlap of $\overline{N}_{\text{cls}}(\kappa)$ and $\overline{N}_{\text{cls}}(\nu)$. Luckily this bad scenario never occurs. We postpone the proof of this property to Lemma 21. The proof of (PD'.3(b)) is similar to that of Lemma 13, and we defer it to Lemma 22.

(SI'.1) follows directly from the algorithm because for each demand $\nu \in j$, all facilities added to $\overline{N}(\nu)$ are immediately removed from $\tilde{N}(j)$ and each facility is added to $\overline{N}(\nu)$ of exactly one demand $\nu \in j$. Splitting facilities obviously preserves (SI'.1).

The proof of (SI'.2) is similar to that of Lemma 11. If $\kappa = \nu$ then (SI'.2) follows from (SI'.1), so we can assume that $\kappa \neq \nu$. Suppose that $\nu' \in j$ is assigned to $\kappa' \in P$ and consider the situation after Phase 1. By the way we reassign facilities in Case 1, at this time we have $\overline{N}(\nu) \subseteq \overline{N}(\kappa) = \overline{N}_{\text{cls}}(\kappa)$ and $\overline{N}(\nu') \subseteq \overline{N}(\kappa') = \overline{N}_{\text{cls}}(\kappa')$, so $\overline{N}(\nu') \cap \overline{N}_{\text{cls}}(\kappa) = \emptyset$, by (PD'.1). Moreover, we have $\tilde{N}(j) \cap \overline{N}_{\text{cls}}(\kappa) = \emptyset$ after this iteration, because any facilities that were also in $\overline{N}_{\text{cls}}(\kappa)$ were removed from $\tilde{N}(j)$ when ν was created. In Phase 2, augmentation does not change $\overline{N}_{\text{cls}}(\kappa)$ and all facilities added to $\overline{N}(\nu')$ are from the set $\tilde{N}(j)$ at the end of Phase 1, which is a subset of the set $\tilde{N}(j)$ after this iteration, since $\tilde{N}(j)$ can only shrink. So the condition (SI'.2) will remain true.

Lemma 21 *Property (PD'.3(a)) holds.*

Proof. Let j be the client for which $\nu \in j$. We consider an iteration when we create ν from j and assign it to κ , and within this proof, notation $\tilde{N}_{\text{cls}}(j)$ and $\tilde{N}(j)$ will refer to the value of the sets at this particular time. At this time, $\overline{N}(\nu)$ is initialized to

$\tilde{N}(j) \cap \overline{N}(\kappa)$. Recall that $\overline{N}(\kappa)$ is now equal to the final $\overline{N}_{\text{cls}}(\kappa)$ (taking into account facility splitting). We would like to show that the set $\tilde{N}_{\text{cls}}(j) \cap \overline{N}_{\text{cls}}(\kappa)$ (which is not empty) will be included in $\overline{N}_{\text{cls}}(\nu)$ at the end. Technically speaking, this will not be true due to facility splitting, so we need to rephrase this claim and the proof in terms of the set of facilities obtained after the algorithm completes.

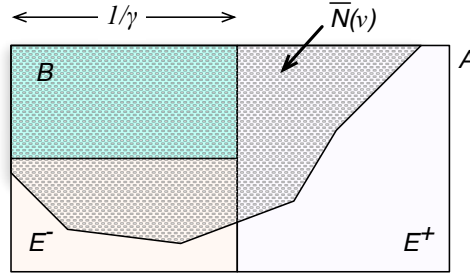


Figure 5.1: Illustration of the sets $\overline{N}(\nu)$, A , B , E^- and E^+ in the proof of Lemma 21. Let $X \subseteq Y$ mean that the facility sets X is obtained from Y by splitting facilities. We then have $A \subseteq \tilde{N}(j)$, $B \subseteq \tilde{N}_{\text{cls}}(j) \cap \overline{N}_{\text{cls}}(\kappa)$, $E^- \subseteq \tilde{N}_{\text{cls}}(j) - \overline{N}_{\text{cls}}(\kappa)$, $E^+ \subseteq \tilde{N}(j) - \tilde{N}_{\text{cls}}(j)$.

We define the sets A , B , E^- and E^+ as the subsets of $\overline{\mathbb{F}}$ (the final set of facilities) that were obtained from splitting facilities in the sets $\tilde{N}(j)$, $\tilde{N}_{\text{cls}}(j) \cap \overline{N}_{\text{cls}}(\kappa)$, $\tilde{N}_{\text{cls}}(j) - \overline{N}_{\text{cls}}(\kappa)$ and $\tilde{N}(j) - \tilde{N}_{\text{cls}}(j)$, respectively. (See Figure 5.1.) We claim that at the end $B \subseteq \overline{N}_{\text{cls}}(\nu)$, with the caveat that the ties in the definition of $\overline{N}_{\text{cls}}(\nu)$ are broken in favor of the facilities in B . (This is the tie-breaking rule that we mentioned in the definition of $\overline{N}_{\text{cls}}(\nu)$.) This will be sufficient to prove the lemma because $B \neq \emptyset$, by the algorithm.

We now prove this claim. In this paragraph $\overline{N}(\nu)$ denotes the final set $\overline{N}(\nu)$ after both phases are completed. Thus the total connection value of $\overline{N}(\nu)$ to ν is 1. Note first that

$B \subseteq \overline{N}(\nu) \subseteq A$, because we never remove facilities from $\overline{N}(\nu)$ and we only add facilities from $\tilde{N}(j)$. Also, $B \cup E^-$ represents the facilities obtained from $\tilde{N}_{\text{cls}}(j)$, so $\sum_{\mu \in B \cup E^-} \bar{y}_\mu = 1/\gamma$. This and $B \subseteq \overline{N}(\nu)$ implies that the total connection value of $B \cup (\overline{N}(\nu) \cap E^-)$ to ν is at most $1/\gamma$. But all facilities in $B \cup (\overline{N}(\nu) \cap E^-)$ are closer to ν (taking into account our tie breaking in property (NB)) than those in $E^+ \cap \overline{N}(\nu)$. It follows that $B \subseteq \overline{N}_{\text{cls}}(\nu)$, completing the proof. ■

Lemma 22 *Property (PD'.3(b)) holds.*

Proof. This proof is similar to that for Lemma 13. For a client j and demand η , we will write $\text{tcc}_{\text{cls}}^\eta(j)$ and $\text{dmax}_{\text{cls}}^\eta(j)$ to denote the values of $\text{tcc}_{\text{cls}}(j)$ and $\text{dmax}_{\text{cls}}(j)$ at the time when η was created. (Here η may or may not be a demand of client j).

Suppose $\nu \in j$ is assigned to a primary demand $\kappa \in p$. By the way primary demands are constructed in the partitioning algorithm, $\tilde{N}_{\text{cls}}(p)$ becomes $\overline{N}(\kappa)$, which is equal to the final value of $\overline{N}_{\text{cls}}(\kappa)$. So we have $C_{\text{cls}}^{\text{avg}}(\kappa) = \text{tcc}_{\text{cls}}^\kappa(p)$ and $C_{\text{cls}}^{\text{max}}(\kappa) = \text{dmax}_{\text{cls}}^\kappa(p)$. Further, since we choose p to minimize $\text{tcc}_{\text{cls}}(p) + \text{dmax}_{\text{cls}}(p)$, we have that $\text{tcc}_{\text{cls}}^\kappa(p) + \text{dmax}_{\text{cls}}^\kappa(p) \leq \text{tcc}_{\text{cls}}^\kappa(j) + \text{dmax}_{\text{cls}}^\kappa(j)$.

Using an argument analogous to that in the proof of Lemma 12, our modified partitioning algorithm guarantees that $\text{tcc}_{\text{cls}}^\kappa(j) \leq \text{tcc}_{\text{cls}}^\nu(j) \leq C_{\text{cls}}^{\text{avg}}(\nu)$ and $\text{dmax}_{\text{cls}}^\kappa(j) \leq \text{dmax}_{\text{cls}}^\nu(j) \leq C_{\text{cls}}^{\text{max}}(\nu)$ since ν was created later. Therefore, we have

$$\begin{aligned} C_{\text{cls}}^{\text{avg}}(\kappa) + C_{\text{cls}}^{\text{max}}(\kappa) &= \text{tcc}_{\text{cls}}^\kappa(p) + \text{dmax}_{\text{cls}}^\kappa(p) \\ &\leq \text{tcc}_{\text{cls}}^\kappa(j) + \text{dmax}_{\text{cls}}^\kappa(j) \leq \text{tcc}_{\text{cls}}^\nu(j) + \text{dmax}_{\text{cls}}^\nu(j) \leq C_{\text{cls}}^{\text{avg}}(\nu) + C_{\text{cls}}^{\text{max}}(\nu), \end{aligned}$$

completing the proof. ■

Now we have completed the proof that the computed partitioning satisfies all the required properties.

Algorithm EBGs. The complete algorithm starts with solving the LP(3.1) and computing the partitioning described earlier in this section. Given the partitioned fractional solution (\bar{x}, \bar{y}) with the desired properties, we start the process of opening facilities and making connections to obtain an integral solution. To this end, for each primary demand $\kappa \in P$, we open exactly one facility $\phi(\kappa)$ in $\bar{N}_{\text{cls}}(\kappa)$, where each $\mu \in \bar{N}_{\text{cls}}(\kappa)$ is chosen as $\phi(\kappa)$ with probability $\gamma \bar{y}_\mu$. For all facilities $\mu \in \bar{\mathbb{F}} - \bigcup_{\kappa \in P} \bar{N}_{\text{cls}}(\kappa)$, we open them independently, each with probability $\gamma \bar{y}_\mu$.

We claim that all probabilities are well-defined, that is $\gamma \bar{y}_\mu \leq 1$ for all μ . Indeed, if $\bar{y}_\mu > 0$ then $\bar{y}_\mu = \bar{x}_{\mu\nu}$ for some ν , by Property (CO). If $\mu \in \bar{N}_{\text{cls}}(\nu)$ then the definition of close neighborhoods implies that $\bar{x}_{\mu\nu} \leq 1/\gamma$. If $\mu \in \bar{N}_{\text{far}}(\nu)$ then $\bar{x}_{\mu\nu} \leq 1 - 1/\gamma \leq 1/\gamma$, because $\gamma < 2$. Thus $\gamma \bar{y}_\mu \leq 1$, as claimed.

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

Analysis. By the algorithm, for each client j , all its r_j demands are connected to open facilities. If two different siblings $\nu, \nu' \in j$ are assigned, respectively, to primary demands κ, κ' then, by Properties (SI'.1), (SI'.2), and (PD'.1) we have

$$(\bar{N}(\nu) \cup \bar{N}_{\text{cls}}(\kappa)) \cap (\bar{N}(\nu') \cup \bar{N}_{\text{cls}}(\kappa')) = \emptyset.$$

This condition guarantees that ν and ν' are assigned to different facilities, regardless whether they are connected to a neighbor facility or to its target facility. Therefore the computed solution is feasible.

We now estimate the cost of the solution computed by Algorithm EBGs. The lemma below bounds the expected facility cost.

Lemma 23 *The expectation of facility cost F_{EBGS} of Algorithm EBGs is at most γF^* .*

Proof. By the algorithm, each facility $\mu \in \bar{\mathbb{F}}$ is opened with probability $\gamma \bar{y}_\mu$, independently of whether it belongs to the close neighborhood of a primary demand or not. Therefore, by linearity of expectation, we have that the expected facility cost is

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

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

In the remainder of this section we focus on the connection cost. Let C_ν be the random variable representing the connection cost of a demand ν . Our objective is to show that the expectation of ν satisfies

$$\mathbb{E}[C_\nu] \leq C^{\text{avg}}(\nu) \cdot \max \left\{ \frac{1/e + 1/e^\gamma}{1 - 1/\gamma}, 1 + \frac{2}{e^\gamma} \right\}. \quad (5.9)$$

If ν is a primary demand then, due to the algorithm, we have $\mathbb{E}[C_\nu] = C_{\text{cls}}^{\text{avg}}(\nu) \leq C^{\text{avg}}(\nu)$, so (5.9) is easily satisfied.

Thus for the rest of the argument we will focus on the case when ν is a non-primary demand. Recall that the algorithm connects ν to the nearest open facility in $\bar{N}_{\text{cls}}(\nu)$ if at least one facility in $\bar{N}_{\text{cls}}(\nu)$ is open. Otherwise the algorithm connects ν to the nearest open

facility in $\overline{N}_{\text{far}}(\nu)$, if any. In the event that no facility in $\overline{N}(\nu)$ opens, the algorithm will connect ν to its target facility $\phi(\kappa)$, where κ is the primary demand that ν was assigned to, and $\phi(\kappa)$ is the only facility open in $\overline{N}_{\text{cls}}(\kappa)$. Let Λ^ν denote the event that at least one facility in $\overline{N}(\nu)$ is open and Λ_{cls}^ν be the event that at least one facility in $\overline{N}_{\text{cls}}(\nu)$ is open. $\neg\Lambda^\nu$ denotes the complement event of Λ^ν , that is, the event that none of ν 's neighbors opens. We want to estimate the following three conditional expectations:

$$\mathbb{E}[C_\nu \mid \Lambda_{\text{cls}}^\nu], \quad \mathbb{E}[C_\nu \mid \Lambda^\nu \wedge \neg\Lambda_{\text{cls}}^\nu], \quad \text{and} \quad \mathbb{E}[C_\nu \mid \neg\Lambda^\nu],$$

and their associated probabilities.

We start with a lemma dealing with the third expectation, $\mathbb{E}[C_\nu \mid \neg\Lambda^\nu] = \mathbb{E}[d_{\phi(\kappa)\nu} \mid \Lambda^\nu]$. The proof of this lemma relies on Properties (PD'.3(a)) and (PD'.3(b)) of modified partitioning and follows the reasoning in the proof of a similar lemma in [6, 5].

Lemma 24 *Assuming that no facility in $\overline{N}(\nu)$ opens, the expected connection cost of ν is*

$$\mathbb{E}[C_\nu \mid \neg\Lambda^\nu] \leq C_{\text{cls}}^{\text{avg}}(\nu) + 2C_{\text{far}}^{\text{avg}}(\nu). \quad (5.10)$$

Proof. It suffices to show a stronger inequality

$$\mathbb{E}[C_\nu \mid \neg\Lambda^\nu] \leq C_{\text{cls}}^{\text{avg}}(\nu) + C_{\text{cls}}^{\text{max}}(\nu) + C_{\text{far}}^{\text{avg}}(\nu), \quad (5.11)$$

which then implies (5.10) because $C_{\text{cls}}^{\text{max}}(\nu) \leq C_{\text{far}}^{\text{avg}}(\nu)$. The proof of (5.11) is similar to that in [5]. For the sake of completeness, we provide it here, formulated in our terminology and notation.

Assume that the event $\neg\Lambda^\nu$ is true, that is Algorithm EBGs does not open any facility in $\overline{N}(\nu)$. Let κ be the primary demand that ν was assigned to. Also let

$$K = \overline{N}_{\text{cls}}(\kappa) \setminus \overline{N}(\nu), \quad V_{\text{cls}} = \overline{N}_{\text{cls}}(\kappa) \cap \overline{N}_{\text{cls}}(\nu) \quad \text{and} \quad V_{\text{far}} = \overline{N}_{\text{cls}}(\kappa) \cap \overline{N}_{\text{far}}(\nu).$$

Then $K, V_{\text{cls}}, V_{\text{far}}$ form a partition of $\overline{N}_{\text{cls}}(\kappa)$, that is, they are disjoint and their union is $\overline{N}_{\text{cls}}(\kappa)$. Moreover, we have that K is not empty, because Algorithm EBGs opens some facility in $\overline{N}_{\text{cls}}(\kappa)$ and this facility cannot be in $V_{\text{cls}} \cup V_{\text{far}}$, by our assumption. We also have that V_{cls} is not empty due to (PD'.3(a)).

Recall that $D(A, \eta) = \sum_{\mu \in A} d_{\mu\eta} \bar{y}_\mu / \sum_{\mu \in A} \bar{y}_\mu$ is the average distance between a demand η and the facilities in a set A . We shall show that

$$D(K, \nu) \leq C_{\text{cls}}^{\text{avg}}(\kappa) + C_{\text{cls}}^{\text{max}}(\kappa) + C_{\text{far}}^{\text{avg}}(\nu). \quad (5.12)$$

This is sufficient, because, by the algorithm, $D(K, \nu)$ is exactly the expected connection cost for demand ν conditioned on the event that none of ν 's neighbors opens, that is the left-hand side of (5.11). Further, (PD'.3(b)) states that $C_{\text{cls}}^{\text{avg}}(\kappa) + C_{\text{cls}}^{\text{max}}(\kappa) \leq C_{\text{cls}}^{\text{avg}}(\nu) + C_{\text{cls}}^{\text{max}}(\nu)$, and thus (5.12) implies (5.11).

The proof of (5.12) is by analysis of several cases.

Case 1: $D(K, \kappa) \leq C_{\text{cls}}^{\text{avg}}(\kappa)$. For any facility $\mu \in V_{\text{cls}}$ (recall that $V_{\text{cls}} \neq \emptyset$), we have $d_{\mu\kappa} \leq C_{\text{cls}}^{\text{max}}(\kappa)$ and $d_{\mu\nu} \leq C_{\text{cls}}^{\text{max}}(\nu) \leq C_{\text{far}}^{\text{avg}}(\nu)$. Therefore, using the case assumption, we get $D(K, \nu) \leq D(K, \kappa) + d_{\mu\kappa} + d_{\mu\nu} \leq C_{\text{cls}}^{\text{avg}}(\kappa) + C_{\text{cls}}^{\text{max}}(\kappa) + C_{\text{far}}^{\text{avg}}(\nu)$.

Case 2: There exists a facility $\mu \in V_{\text{cls}}$ such that $d_{\mu\kappa} \leq C_{\text{cls}}^{\text{avg}}(\kappa)$. Since $\mu \in V_{\text{cls}}$, we infer that $d_{\mu\nu} \leq C_{\text{cls}}^{\text{max}}(\nu) \leq C_{\text{far}}^{\text{avg}}(\nu)$. Using $C_{\text{cls}}^{\text{max}}(\kappa)$ to bound $D(K, \kappa)$, we have $D(K, \nu) \leq D(K, \kappa) + d_{\mu\kappa} + d_{\mu\nu} \leq C_{\text{cls}}^{\text{max}}(\kappa) + C_{\text{cls}}^{\text{avg}}(\kappa) + C_{\text{far}}^{\text{avg}}(\nu)$.

Case 3: In this case we assume that neither of Cases 1 and 2 applies, that is $D(K, \kappa) > C_{\text{cls}}^{\text{avg}}(\kappa)$ and every $\mu \in V_{\text{cls}}$ satisfies $d_{\mu\kappa} > C_{\text{cls}}^{\text{avg}}(\kappa)$. This implies that $D(K \cup V_{\text{cls}}, \kappa) > C_{\text{cls}}^{\text{avg}}(\kappa) = D(\overline{N}_{\text{cls}}(\kappa), \kappa)$. Since sets K, V_{cls} and V_{far} form a partition of $\overline{N}_{\text{cls}}(\kappa)$, we obtain

that in this case V_{far} is not empty and $D(V_{\text{far}}, \kappa) < C_{\text{cls}}^{\text{avg}}(\kappa)$. Let $\delta = C_{\text{cls}}^{\text{avg}}(\kappa) - D(V_{\text{far}}, \kappa) > 0$.

We now have two sub-cases:

Case 3.1: $D(V_{\text{far}}, \nu) \leq C_{\text{far}}^{\text{avg}}(\nu) + \delta$. Substituting δ , this implies that $D(V_{\text{far}}, \nu) + D(V_{\text{far}}, \kappa) \leq$

$C_{\text{cls}}^{\text{avg}}(\kappa) + C_{\text{far}}^{\text{avg}}(\nu)$. From the definition of the average distance $D(V_{\text{far}}, \kappa)$ and $D(V_{\text{far}}, \nu)$,

we obtain that there exists some $\mu \in V_{\text{far}}$ such that $d_{\mu\kappa} + d_{\mu\nu} \leq C_{\text{cls}}^{\text{avg}}(\kappa) + C_{\text{far}}^{\text{avg}}(\nu)$.

Thus $D(K, \nu) \leq D(K, \kappa) + d_{\mu\kappa} + d_{\mu\nu} \leq C_{\text{cls}}^{\text{max}}(\kappa) + C_{\text{cls}}^{\text{avg}}(\kappa) + C_{\text{far}}^{\text{avg}}(\nu)$.

Case 3.2: $D(V_{\text{far}}, \nu) > C_{\text{far}}^{\text{avg}}(\nu) + \delta$. The case assumption implies that V_{far} is a proper

subset of $\bar{N}_{\text{far}}(\nu)$, that is $\bar{N}_{\text{far}}(\nu) \setminus V_{\text{far}} \neq \emptyset$. Let $\hat{y} = \gamma \sum_{\mu \in V_{\text{far}}} \bar{y}_{\mu}$. We can express

$C_{\text{far}}^{\text{avg}}(\nu)$ using \hat{y} as follows

$$C_{\text{far}}^{\text{avg}}(\nu) = D(V_{\text{far}}, \nu) \frac{\hat{y}}{\gamma - 1} + D(\bar{N}_{\text{far}}(\nu) \setminus V_{\text{far}}, \nu) \frac{\gamma - 1 - \hat{y}}{\gamma - 1}.$$

Then, using the case condition and simple algebra, we have

$$\begin{aligned} C_{\text{cls}}^{\text{max}}(\nu) &\leq D(\bar{N}_{\text{far}}(\nu) \setminus V_{\text{far}}, \nu) \\ &\leq C_{\text{far}}^{\text{avg}}(\nu) - \frac{\hat{y}\delta}{\gamma - 1 - \hat{y}} \leq C_{\text{far}}^{\text{avg}}(\nu) - \frac{\hat{y}\delta}{1 - \hat{y}}, \end{aligned} \quad (5.13)$$

where the last step follows from $1 < \gamma < 2$.

On the other hand, since K , V_{cls} , and V_{far} form a partition of $\bar{N}_{\text{cls}}(\kappa)$, we have

$C_{\text{cls}}^{\text{avg}}(\kappa) = (1 - \hat{y})D(K \cup V_{\text{cls}}, \kappa) + \hat{y}D(V_{\text{far}}, \kappa)$. Then using the definition of δ we obtain

$$D(K \cup V_{\text{cls}}, \kappa) = C_{\text{cls}}^{\text{avg}}(\kappa) + \frac{\hat{y}\delta}{1 - \hat{y}}. \quad (5.14)$$

Now we are essentially done. If there exists some $\mu \in V_{\text{cls}}$ such that $d_{\mu\kappa} \leq C_{\text{cls}}^{\text{avg}}(\kappa) +$

$\hat{y}\delta/(1 - \hat{y})$, then we have

$$\begin{aligned}
D(K, \nu) &\leq D(K, \kappa) + d_{\mu\kappa} + d_{\mu\nu} \\
&\leq C_{\text{cls}}^{\text{max}}(\kappa) + C_{\text{cls}}^{\text{avg}}(\kappa) + \frac{\hat{y}\delta}{1 - \hat{y}} + C_{\text{cls}}^{\text{max}}(\nu) \\
&\leq C_{\text{cls}}^{\text{max}}(\kappa) + C_{\text{cls}}^{\text{avg}}(\kappa) + C_{\text{far}}^{\text{avg}}(\nu),
\end{aligned}$$

where we used (5.13) in the last step. Otherwise, from (5.14), we must have $D(K, \kappa) \leq$

$C_{\text{cls}}^{\text{avg}}(\kappa) + \hat{y}\delta/(1 - \hat{y})$. Choosing any $\mu \in V_{\text{cls}}$, it follows that

$$\begin{aligned}
D(K, \nu) &\leq D(K, \kappa) + d_{\mu\kappa} + d_{\mu\nu} \\
&\leq C_{\text{cls}}^{\text{avg}}(\kappa) + \frac{\hat{y}\delta}{1 - \hat{y}} + C_{\text{cls}}^{\text{max}}(\kappa) + C_{\text{cls}}^{\text{max}}(\nu) \\
&\leq C_{\text{cls}}^{\text{avg}}(\kappa) + C_{\text{cls}}^{\text{max}}(\kappa) + C_{\text{far}}^{\text{avg}}(\nu),
\end{aligned}$$

again using (5.13) in the last step.

This concludes the proof of (5.10). As explained earlier, Lemma 24 follows. ■

Next, we derive some estimates for the expected cost of direct connections. The next technical lemma is a generalization of Lemma 18. In Lemma 18 we bound the expected distance to the closest open facility in $\overline{N}(\nu)$, conditioned on at least one facility in $\overline{N}(\nu)$ being open. The lemma below provides a similar estimate for an arbitrary set A of facilities in $\overline{N}(\nu)$, conditioned on that at least one facility in set A is open. Recall that $D(A, \nu) = \sum_{\mu \in A} d_{\mu\nu} \bar{y}_\mu / \sum_{\mu \in A} \bar{y}_\mu$ is the average distance from ν to a facility in A .

Lemma 25 *For any non-empty set $A \subseteq \overline{N}(\nu)$, let Λ_A^ν be the event that at least one facility in A is opened by Algorithm EBGs, and denote by $C_\nu(A)$ the random variable representing*

the distance from ν to the closest open facility in A . Then the expected distance from ν to the nearest open facility in A , conditioned on at least one facility in A being opened, is

$$\mathbb{E}[C_\nu(A) \mid \Lambda_A^\nu] \leq D(A, \nu).$$

Proof. The proof follows the same reasoning as the proof of Lemma 18, so we only sketch it here. We start with a similar grouping of facilities in A : for each primary demand κ , if $\overline{N}_{\text{cls}}(\kappa) \cap A \neq \emptyset$ then $\overline{N}_{\text{cls}}(\kappa) \cap A$ forms a group. Facilities in A that are not in a neighborhood of any primary demand form singleton groups. We denote these groups G_1, \dots, G_k . It is clear that the groups are disjoint because of (PD'.1). Denoting by $\bar{d}_s = D(G_s, \nu)$ the average distance from ν to a group G_s , we can assume that these groups are ordered so that $\bar{d}_1 \leq \dots \leq \bar{d}_k$.

Each group can have at most one facility open and the events representing opening of any two facilities that belong to different groups are independent. To estimate the distance from ν to the nearest open facility in A , we use an alternative random process to make connections, that is easier to analyze. Instead of connecting ν to the nearest open facility in A , we will choose the smallest s for which G_s has an open facility and connect ν to this facility. (Thus we selected an open facility with respect to the minimum \bar{d}_s , not the actual distance from ν to this facility.) This can only increase the expected connection cost, thus denoting $g_s = \sum_{\mu \in G_s} \gamma \bar{y}_\mu$ for all $s = 1, \dots, k$, and letting $\mathbb{P}[\Lambda_A^\nu]$ be the probability that

A has at least one facility open, we have

$$\mathbb{E}[C_\nu(A) \mid \Lambda_A^\nu] \leq \frac{1}{\mathbb{P}[\Lambda_A^\nu]} (\bar{d}_1 g_1 + \bar{d}_2 g_2 (1 - g_1) + \dots + \bar{d}_k g_k (1 - g_1) \dots (1 - g_{k-1})) \quad (5.15)$$

$$\begin{aligned} &\leq \frac{1}{\mathbb{P}[\Lambda_A^\nu]} \frac{\sum_{s=1}^k \bar{d}_s g_s}{\sum_{s=1}^k g_s} (1 - \prod_{s=1}^k (1 - g_s)) \\ &= \frac{\sum_{s=1}^k \bar{d}_s g_s}{\sum_{s=1}^k g_s} = \frac{\sum_{\mu \in A} d_{\mu\nu} \gamma \bar{y}_\mu}{\sum_{\mu \in A} \gamma \bar{y}_\mu} \\ &= \frac{\sum_{s=1}^k d_{\mu\nu} \bar{y}_\mu}{\sum_{\mu \in A} \bar{y}_\mu} = D(A, \nu). \end{aligned} \quad (5.16)$$

Inequality (5.16) follows from inequality (A.3) in A.2. The rest of the derivation follows from $\mathbb{P}[\Lambda_A^\nu] = 1 - \prod_{s=1}^k (1 - g_s)$, and the definition of \bar{d}_s , g_s and $D(A, \nu)$. ■

A consequence of Lemma 25 is the following corollary which bounds the other two expectations of C_ν , when at least one facility is opened in $\bar{N}_{\text{cls}}(\nu)$, and when no facility in $\bar{N}_{\text{cls}}(\nu)$ opens but a facility in $\bar{N}_{\text{far}}(\nu)$ is opened.

Corollary 26 (a) $\mathbb{E}[C_\nu \mid \Lambda_{\text{cls}}^\nu] \leq C_{\text{cls}}^{\text{avg}}(\nu)$, and (b) $\mathbb{E}[C_\nu \mid \Lambda^\nu \wedge \neg \Lambda_{\text{cls}}^\nu] \leq C_{\text{far}}^{\text{avg}}(\nu)$.

Proof. When there is an open facility in $\bar{N}_{\text{cls}}(\nu)$, the algorithm connect ν to the nearest open facility in $\bar{N}_{\text{cls}}(\nu)$. When no facility in $\bar{N}_{\text{cls}}(\nu)$ opens but some facility in $\bar{N}_{\text{far}}(\nu)$ opens, the algorithm connects ν to the nearest open facility in $\bar{N}_{\text{far}}(\nu)$. The rest of the proof follows from Lemma 25. By setting the set A in Lemma 25 to $\bar{N}_{\text{cls}}(\nu)$, we have

$$\mathbb{E}[C_\nu \mid \Lambda_{\text{cls}}^\nu] \leq D(\bar{N}_{\text{cls}}(\nu), \nu) = C_{\text{cls}}^{\text{avg}}(\nu),$$

proving part (a), and by setting the set A to $\bar{N}_{\text{far}}(\nu)$, we have

$$\mathbb{E}[C_\nu \mid \Lambda^\nu \wedge \neg \Lambda_{\text{cls}}^\nu] \leq D(\bar{N}_{\text{far}}(\nu), \nu) = C_{\text{far}}^{\text{avg}}(\nu),$$

which proves part (b). ■

Given the estimate on the three expected distances when ν connects to its close facility in $\overline{N}_{\text{cls}}(\nu)$ in (5.3), or its far facility in $\overline{N}_{\text{far}}(\nu)$ in (5.3), or its target facility $\phi(\kappa)$ in (5.10), the only missing pieces are estimates on the corresponding probabilities of each event, which we do in the next lemma. Once done, we shall put all pieces together and proving the desired inequality on $\mathbb{E}[C_\nu]$, that is (5.9).

The next Lemma bounds the probabilities for events that no facilities in $\overline{N}_{\text{cls}}(\nu)$ and $\overline{N}(\nu)$ are opened by the algorithm.

Lemma 27 (a) $\mathbb{P}[\neg\Lambda_{\text{cls}}^\nu] \leq 1/e$, and (b) $\mathbb{P}[\neg\Lambda^\nu] \leq 1/e^\gamma$.

Proof. (a) To estimate $\mathbb{P}[\neg\Lambda_{\text{cls}}^\nu]$, we again consider a grouping of facilities in $\overline{N}_{\text{cls}}(\nu)$, as in the proof of Lemma 25, according to the primary demand's close neighborhood that they fall in, with facilities not belonging to such neighborhoods forming their own singleton groups. As before, the groups are denoted G_1, \dots, G_k . It is easy to see that $\sum_{s=1}^k g_s = \sum_{\mu \in \overline{N}_{\text{cls}}(\nu)} \gamma \bar{y}_\mu = 1$. For any group G_s , the probability that a facility in this group opens is $\sum_{\mu \in G_s} \gamma \bar{y}_\mu = g_s$ because in the algorithm at most one facility in a group can be chosen and each is chosen with probability $\gamma \bar{y}_\mu$. Therefore the probability that no facility opens is $\prod_{s=1}^k (1 - g_s)$, which is at most $e^{-\sum_{s=1}^k g_s} = 1/e$. Therefore we have $\mathbb{P}[\neg\Lambda_A^\nu] \leq 1/e$.

(b) This proof is similar to the proof of (a). The probability $\mathbb{P}[\neg\Lambda^\nu]$ is at most $e^{-\sum_{s=1}^k g_s} = 1/e^\gamma$, because we now have $\sum_{s=1}^k g_s = \gamma \sum_{\mu \in \overline{N}(\nu)} \bar{y}_\mu = \gamma \cdot 1 = \gamma$. ■

We are now ready to bound the overall connection cost of Algorithm EBGs, namely inequality (5.9).

Lemma 28 *The expected connection of ν is*

$$\mathbb{E}[C_\nu] \leq C^{\text{avg}}(\nu) \cdot \max \left\{ \frac{1/e + 1/e^\gamma}{1 - 1/\gamma}, 1 + \frac{2}{e^\gamma} \right\}.$$

Proof. Recall that, to connect ν , the algorithm uses the closest facility in $\overline{N}_{\text{cls}}(\nu)$ if one is opened; otherwise it will try to connect ν to the closest facility in $\overline{N}_{\text{far}}(\nu)$. Failing that, it will connect ν to $\phi(\kappa)$, the sole facility open in the neighborhood of κ , the primary demand ν was assigned to. Given that, we estimate $\mathbb{E}[C_\nu]$ as follows:

$$\begin{aligned} \mathbb{E}[C_\nu] &= \mathbb{E}[C_\nu \mid \Lambda_{\text{cls}}^\nu] \cdot \mathbb{P}[\Lambda_{\text{cls}}^\nu] + \mathbb{E}[C_\nu \mid \Lambda^\nu \wedge \neg \Lambda_{\text{cls}}^\nu] \cdot \mathbb{P}[\Lambda^\nu \wedge \neg \Lambda_{\text{cls}}^\nu] \\ &\quad + \mathbb{E}[C_\nu \mid \neg \Lambda^\nu] \cdot \mathbb{P}[\neg \Lambda^\nu] \\ &\leq C_{\text{cls}}^{\text{avg}}(\nu) \cdot \mathbb{P}[\Lambda_{\text{cls}}^\nu] + C_{\text{far}}^{\text{avg}}(\nu) \cdot \mathbb{P}[\Lambda^\nu \wedge \neg \Lambda_{\text{cls}}^\nu] \\ &\quad + [C_{\text{cls}}^{\text{avg}}(\nu) + 2C_{\text{far}}^{\text{avg}}(\nu)] \cdot \mathbb{P}[\neg \Lambda^\nu] \\ &= [C_{\text{cls}}^{\text{avg}}(\nu) + C_{\text{far}}^{\text{avg}}(\nu)] \cdot \mathbb{P}[\neg \Lambda^\nu] + [C_{\text{far}}^{\text{avg}}(\nu) - C_{\text{cls}}^{\text{avg}}(\nu)] \cdot \mathbb{P}[\neg \Lambda_{\text{cls}}^\nu] + C_{\text{cls}}^{\text{avg}}(\nu) \\ &\leq [C_{\text{cls}}^{\text{avg}}(\nu) + C_{\text{far}}^{\text{avg}}(\nu)] \cdot \frac{1}{e^\gamma} + [C_{\text{far}}^{\text{avg}}(\nu) - C_{\text{cls}}^{\text{avg}}(\nu)] \cdot \frac{1}{e} + C_{\text{cls}}^{\text{avg}}(\nu) \\ &= \left(1 - \frac{1}{e} + \frac{1}{e^\gamma}\right) \cdot C_{\text{cls}}^{\text{avg}}(\nu) + \left(\frac{1}{e} + \frac{1}{e^\gamma}\right) \cdot C_{\text{far}}^{\text{avg}}(\nu). \end{aligned} \tag{5.18}$$

Inequality (5.17) follows from Corollary 26 and Lemma 24. Inequality (5.18) follows from Lemma 27 and $C_{\text{far}}^{\text{avg}}(\nu) - C_{\text{cls}}^{\text{avg}}(\nu) \geq 0$.

Now define $\rho = C_{\text{cls}}^{\text{avg}}(\nu)/C^{\text{avg}}(\nu)$. It is easy to see that ρ is between 0 and 1.

Continuing the above derivation, applying (5.8), we get

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

and the proof is now complete. ■

With Lemma 28 proven, we are now ready to bound our total connection cost.

For any client j we have

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

Summing over all clients j we obtain that the total expected connection cost is

$$\mathbb{E}[C_{\text{EBGS}}] \leq C^* \max \left\{ \frac{1/e + 1/e^\gamma}{1 - 1/\gamma}, 1 + \frac{2}{e^\gamma} \right\}.$$

Recall that the expected facility cost is bounded by γF^* , as argued earlier. Hence the total expected 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 29 *Algorithm EBGs is a 1.575-approximation algorithm for FTFP.*

Chapter 6

Primal-dual Algorithms

In this chapter we present results and discuss combinatorial algorithms to the FTFP problem. These approaches, although employ the Linear Program in guiding the algorithm and deriving approximation ratio, the use of LP is implicit. In particular, the algorithms do not require solving the LP and having access to a fractional optimal solution. Two notable such approaches are primal-dual and dual-fitting. In this chapter we assume the primal problem is a minimization problem and the dual problem is a maximization problem, to be consistent with our FTFP problem. In primal-dual algorithms, we start with a feasible dual solution, usually with all dual variables set to zero, then we raise a subset of dual variables and update the corresponding primal variables accordingly. At any time, we keep the dual solution feasible and we stop when the primal solution becomes feasible. The approximation ratio is derived by a relaxed version of the complementary slackness conditions.

Another approach, dual-fitting, starts with an empty dual solution as well, and

works in iterations. In each iteration, the algorithm raises a subset of dual variables, and updates corresponding primal variables. The algorithm stops when the primal solution is feasible. The difference is that in dual-fitting, the dual solution may not be feasible, and we require the cost of the primal solution bounded by the value of the possibly infeasible dual solution. The next step is to find a suitable number γ , which may depend on the input size, such that the dual solution, when divided by γ , becomes feasible. It is easy to see that γ provides an upper bound on the approximation ratio, because the value of a feasible dual solution is a lower bound on the value of an optimal primal solution.

Jain and Vazirani [20] designed a primal-dual algorithm, which we call the JV algorithm, for the UFL problem. Recall that in the UFL, all demands $r_j = 1$. In the JV algorithm, every client j has a number α_j associated. All α_j start at zero, and all clients are unconnected initially. The algorithm has two phases. The first phase runs in iterations. In each iteration, all α_j that were not temporarily connected are raised uniformly. The contribution from a client j to a facility i is $\max\{0, \alpha_j - d_{ij}\}$. Whenever a facility received enough total contribution, that is $\sum_{j \in \mathbb{C}} (\alpha_j - d_{ij})_+ = f_i$, then i is temporarily open and all clients with $\alpha_j \geq d_{ij}$ temporarily connect to i . The facility i is called the witness of the client j . The first phase concludes when all clients are temporarily connected. In the second phase, we construct an auxiliary graph with nodes being temporarily open facilities in the first phase. Two nodes i_1 and i_2 are connected by an edge in this auxiliary graph if there exists some client j that contributes to both of them, that is, $\alpha_j > d_{i_1 j}$ and $\alpha_j > d_{i_2 j}$. We then pick a maximal independent set in the auxiliary graph as the set of facilities to open. For connections, if a client j has an open facility i with $d_{ij} < \alpha_j$, then it connects to that

facility. Notice that there can be at most one such facility. If not, and if there exists some open facility i such that $d_{ij} = \alpha_j$, then j connects to i . These two types of connections are called *direct* connections. If neither is true, then there exists some temporarily open facility i such that $\alpha_j \geq d_{ij}$, namely j 's witness. Since i is not open, there must exist some facility i' that is open and some client j' such that $\alpha_{j'} > d_{i'j}$ and $\alpha_{j'} > d_{ij'}$. We then connect j to one such i' and this type of connection is an *indirect* connection.

We now analyze the cost of the solution. Due to the algorithm, we have $\alpha_{j'} \leq \min\{t(i), t(i')\}$ where $t(i)$ is the time that facility i is temporarily open. The reason is that, if $\alpha_{j'} > t(i)$, then it would have temporarily connected to facility i earlier so its $\alpha_{j'}$ value would have been smaller. On the other hand, since facility i is the witness of client j , we have $t(i) \leq \alpha_j$ ¹. Therefore we have $\alpha_j \geq \alpha_{j'} \geq \max\{d_{i'j'}, d_{ij'}\}$. In addition, we also have $\alpha_j \geq d_{ij}$. Hence $d_{ij} \leq d_{i'j'} + d_{ij'} + d_{ij} \leq \alpha_{j'} + \alpha_{j'} + \alpha_j \leq 3\alpha_j$. To complete the dual solution, we define $\beta_{ij} = \alpha_j - d_{ij}$ if facility i is open and client j contributes to facility i , and $\beta_{ij} = 0$ otherwise. Now we estimate the total cost of this dual solution. For facility cost we have $\sum_{i,j} \beta_{ij}$, and for connection cost, if a client j is directly connected, then its $d_{ij} \leq \alpha_j - \beta_{ij}$, otherwise it is $d_{ij} \leq 3\alpha_j$. The total cost is hence no more than $3 \sum_{j \in \mathbb{C}} \alpha_j$. Since $\{\alpha_j\}$ form a feasible dual solution, we have the optimal solution value is no less than $\sum_{j \in \mathbb{C}} \alpha_j$. Therefore, we have our solution costs no more than 3 times the cost of an optimal solution.

The fault-tolerant facility location problem (FTFL) was introduced by Jain and Vazirani [21] primarily to demonstrate that their primal-dual algorithm on UFL can be

¹ $t(i) < \alpha_j$ is possible if facility i is temporarily open and later j has $\alpha_j = d_{ij}$ to temporarily connect to facility i .

applied to a more general problem, where clients could have demand more than 1, and each facility could be open or close. A client j with demand r_j needs to be connected to r_j different facilities. The primal-dual algorithm by Jain and Vazirani on FTFL gives a ratio of $3 \ln \max_j r_j$. Subsequent attempts on adapting either the primal-dual approach or the dual-fitting approach to FTFL with a sub-logarithmic approximation ratio were not successful, although for the uniform-demand case, that is, when all r_j are equal, Adrian Bumb [3] showed that the JV algorithm [20] for UFL can be adapted to obtain the same ratio 3 as for UFL. On a separate paper, Swamy and Shmoys [27] showed that a greedy algorithm analyzed using dual-fitting can be shown to have a ratio of 1.52. For the non-uniform demand case, the best known result is an $O(\log n)$ -approximation [21].

For our problem, FTFP, we have seen in earlier chapters that it can be approximated with the same ratio as UFL when LP-rounding is used. However, the attempt to obtain a sub-logarithmic approximation ratio on FTFP using the primal-dual algorithm or the dual-fitting algorithm were not successful. On the positive side, we derive a weak result that the greedy algorithm does give a $O(\log n)$ ratio for FTFP. This also sets stage for our presentation of a hard example in the following section. We remark here that the $O(\log n)$ ratio does not even use the triangle inequality. On the negative side, we provide an example showing that the greedy algorithm with dual-fitting analysis can at best give a ratio of $O(\log n / \log \log n)$ under a very reasonable assumption, which we call the *local-charging* assumption. As usual we have $n = |\mathcal{C}|$.

In the following we first describe the greedy algorithm and its analysis. We show that the greedy algorithm is $O(\log n)$ -approximation using the dual-fitting analysis. Then

we present our example showing a lower bound on the dual-fitting analysis on the greedy algorithm. We conclude this chapter with some possible approaches to obtain sub-logarithmic approximation results.

6.1 The Greedy algorithm with $O(\log n)$ Ratio

In this section we show that the greedy algorithm which repeatedly picking the best star (the one with minimum average cost) gives an approximation ratio of $H_n = \ln(n)$ where $n = |\mathbb{C}|$ is the number of clients. A star is a site i and a subset of clients C' . The cost of such a star S is $c(S) = f_i + \sum_{j \in C'} d_{ij}$, and the average cost of S is $c(S)/|C'|$. Call a client j fully-connected, or *exhausted* if j has made r_j connections. Let U be the set of not fully-connected clients. While not all clients fully-connected, the algorithm picks a star $S = (i, C')$ with $C' \subseteq U$, and open one facility at site i . Each client in C' then makes one more connection with site i . The algorithm terminates when all clients are fully-connected. To see the algorithm can be implemented to run in polynomial time, one observes that once a star becomes best, it remains best until one or more of its member clients become exhausted. Thus we can accomplish multiple iterations in a single step and the number of steps is polynomially bounded by $|\mathbb{F}| \cdot |\mathbb{C}|$.

Now we analyze the cost of the solution, using the dual-fitting analysis similar to Jain *et al.* [18]. When we run the greedy algorithm, for every client j , we associate each demand of j with a number α_j^l , which is the average cost of the star when l^{th} demand of j is connected. Now we let $\alpha_j = \alpha_j^{r_j}$, that is, take α_j to be the finishing α_j^l , and order clients

by increasing α_j . That is,

$$\alpha_1 \leq \alpha_2 \leq \dots \leq \alpha_n$$

Due to the algorithm, for every $j = 1, \dots, n$, we have

$$\sum_{l=j}^n (\alpha_j - d_{il})_+ \leq f_i$$

for every site i . The reason is that, when the last demand of j is connected, all clients $j+1, \dots, n$ are still active so (their α_j are still increasing) their total contribution cannot exceed f_i .

Now we take a closer look at the numbers $\{\alpha_j\}$. We know that the algorithm's total cost is exactly $\sum_{j=1}^n \sum_{l=1}^{r_j} \alpha_j^l$, which is no more than $\sum_{j=1}^n r_j \alpha_j$ since we take α_j to be $\alpha_j^{r_j}$. Now if we can show that $\sum_{j=1}^n r_j \alpha_j$ is no more than $\gamma \cdot \text{OPT}$, where OPT is the cost of an integral optimal solution to the FTFP instance, then we claim our algorithm computes an integral solution whose cost is within a factor of γ from OPT.

We show that $\sum_{j=1}^n r_j \alpha_j$ is within a factor of γ from OPT by showing that $\{\alpha_j/\gamma\}$ is a feasible dual solution to the following program, which is the dual program of the primal LP for FTFP.

$$\begin{aligned} & \max \sum_j r_j \alpha_j \\ & \text{subject to: } \sum_{j=1}^n (\alpha_j - d_{ij})_+ \leq f_i \quad \forall i \in \mathbb{F} \end{aligned}$$

To find the minimum γ that would shrink $\{\alpha_j\}$ to a feasible dual solution, we need to find a worst case instance to maximize γ , also it is clear that the worst case instance

must contain a star whose feasibility requirement would achieves the value of γ , and this star would be the worst star in that instance.

As a first step we can drop the $\max\{0, \cdot\}$, because we can always find a new star by dropping those j with $\alpha_j - d_{ij}$ term negative, and that new star would still be a worst case star. Suppose a worst case star has k clients, and is with facility i , then we have

$$\sum_{j=1}^k \alpha_j - d_{ij} \leq f_i$$

Here we rename clients in the new star to be $1, \dots, k$, although among them, they are still ordered by their α_j .

Now our goal is to find a supremum of the following program:

$$\begin{aligned} & \max \frac{\sum_{j=1}^k \alpha_j}{f_i + \sum_{j=1}^k d_{ij}} \\ & \text{subject to: } \sum_{l=j}^k (\alpha_l - d_{il})_+ \leq f_i \text{ for } j = 1, \dots, k \end{aligned}$$

Actually this is a series of programs for $k = 1, 2, \dots$

Since we are dealing with a particular star, we can abstract away i , to obtain the following program:

$$\begin{aligned} & \max \frac{\sum_{j=1}^k \alpha_j}{f + \sum_{j=1}^k d_j} \\ & \text{subject to: } \sum_{l=j}^k (\alpha_j - d_l)_+ \leq f \text{ for } j = 1, \dots, k \end{aligned} \tag{6.1}$$

Now we claim we can drop the $\max\{0, \cdot\}$ operator because this would relax the constraint in (6.1) and can only make objective value larger (since we are maximizing), so the real optimal is upper bounded by the relaxed optimal. This allows us to work on the following program instead.

$$\begin{aligned} & \max \frac{\sum_{j=1}^k \alpha_j}{f + \sum_{j=1}^k d_j} \\ & \text{subject to: } \sum_{l=j}^k (\alpha_l - d_l) \leq f \text{ for } j = 1, \dots, k \end{aligned} \tag{6.2}$$

For each $j = 1, \dots, k$, the constraint above simply can be rewritten as

$$(k - j + 1)\alpha_j \leq f + \sum_{l=j}^k d_l \leq f + \sum_{l=1}^k d_l. \tag{6.3}$$

The first inequality is a rewrite of the constraint in (6.2) and the second is straightforward.

Therefore we have $\alpha_j \leq (1/(k - j + 1))(f + \sum_{l=j}^k d_l)$, and it easily follows that

$$\sum_{j=1}^n \alpha_j \leq (1/k + 1/(k-1) + \dots + 1) = H_k \leq H_n = \ln(n) \tag{6.4}$$

So we have shown that $\{\alpha_j/\gamma\}$ when $\gamma = H_n$ is a feasible dual solution and therefore $\sum_{j \in \mathbb{C}} r_j \alpha_j$ is no more than $H_n \cdot \text{OPT}$. As our primal solution has cost no more than $\sum_{j \in \mathbb{C}} r_j \alpha_j$, the greedy algorithm computes a solution within H_n from the optimal.

The H_n approximation result is rather weak and is hardly the best approximation ratio possible for the greedy algorithm. An astute reader might notice that we have not used the triangle inequality in deriving the approximation ratio, although we are working on metric FTFP. On the other hand, similar attempts in adapting the greedy algorithm for

UFL to FTFL and look for a sub-logarithmic ratio were not successful by other researchers, as described in the beginning of this chapter. Although FTFP seems to be easier to approximate than FTFL when LP-rounding algorithms were used, it seems the fault-tolerant requirement in both problems sets a hurdle for primal-dual based techniques. In the following section we provide an example that illustrates some difficulty when adapting the dual-fitting analysis to fault-tolerant facility location problems.

6.2 An Example Showing the Difficulty in Obtaining $O(1)$ Ratio

For FTFP, the greedy algorithm that repeatedly picks the best star until all clients have all demands satisfied can be implemented in polynomial time. In Section 6.1 we show that this algorithm is H_n -approximation where $n = |\mathcal{C}|$ is the number of clients. Since the same greedy algorithm is shown to have constant approximation ratio for UFL [23], a natural question to ask is whether greedy can be shown to have $O(1)$ approximation ratio. Here we give an example that hints a negative answer.

We assume the greedy algorithm is analyzed using the dual-fitting technique, which associates with every client j with a number α_j , interpreted as a dual solution to the LP (3.2). However, the dual solution $\{\alpha_j\}$ in general may not be feasible. The dual-fitting technique aims at finding a smallest possible number γ such that, after the dual solution $\{\alpha_j\}$ is shrunk (divided) by γ , all dual constraints are satisfied. That is

$$\sum_{j \in \mathcal{C}} (\alpha_j / \gamma - d_{ij})_+ \leq f_i \quad \text{for all } i \in \mathcal{F}.$$

That γ is taken as the approximation ratio.

In the greedy algorithm, a star with minimum average cost is picked at each iteration and each member client of that star then gets one more connection. It is not specified by the algorithm how we distribute the cost of f_i into member clients, which is part of the analysis. Nonetheless we assume that the cost of f_i is distributed among members only, and not to clients outside this star. We call this *local charging* assumption. Our second assumption is that the proposed dual solution α_j , is taken as the average of individual α_j^l for each of the l^{th} demand of client j , with $l = 1, \dots, r_j$. That is $\alpha_j = \sum_{l=1}^{r_j} \alpha_j^l / r_j$. Suppose the l^{th} demand of j is satisfied while j is in a star with facility i , then $\alpha_j^l = d_{ij} + f_i^{j,l}$, where $f_i^{j,l}$ is the portion of f_i attributed to j in the analysis. Notice that taking the average implies the α_j values thus computed make $\sum_{j \in C} r_j \alpha_j$ equal to the cost of the integral solution by the greedy algorithm.

We now give our example in Figure 6.1. Our example has one site and k groups of clients. Opening one facility at that site costs f_1 . The first group has n_1 clients each with demand r_1 , all at distance $d_1 = 0$ from f_1 . The other groups are listed below:

$$d_1 = 0$$

$$d_2 = \frac{f_1}{n_1}$$

$$d_3 = f_1/n_2 + d_2 = f_1/n_2 + f_1/n_1 = f_1\left(\frac{1}{n_2} + \frac{1}{n_1}\right)$$

...

$$d_k = f_1/n_{k-1} + d_{k-1} = f_1\left(\frac{1}{n_{k-1}} + \dots + \frac{1}{n_1}\right)$$

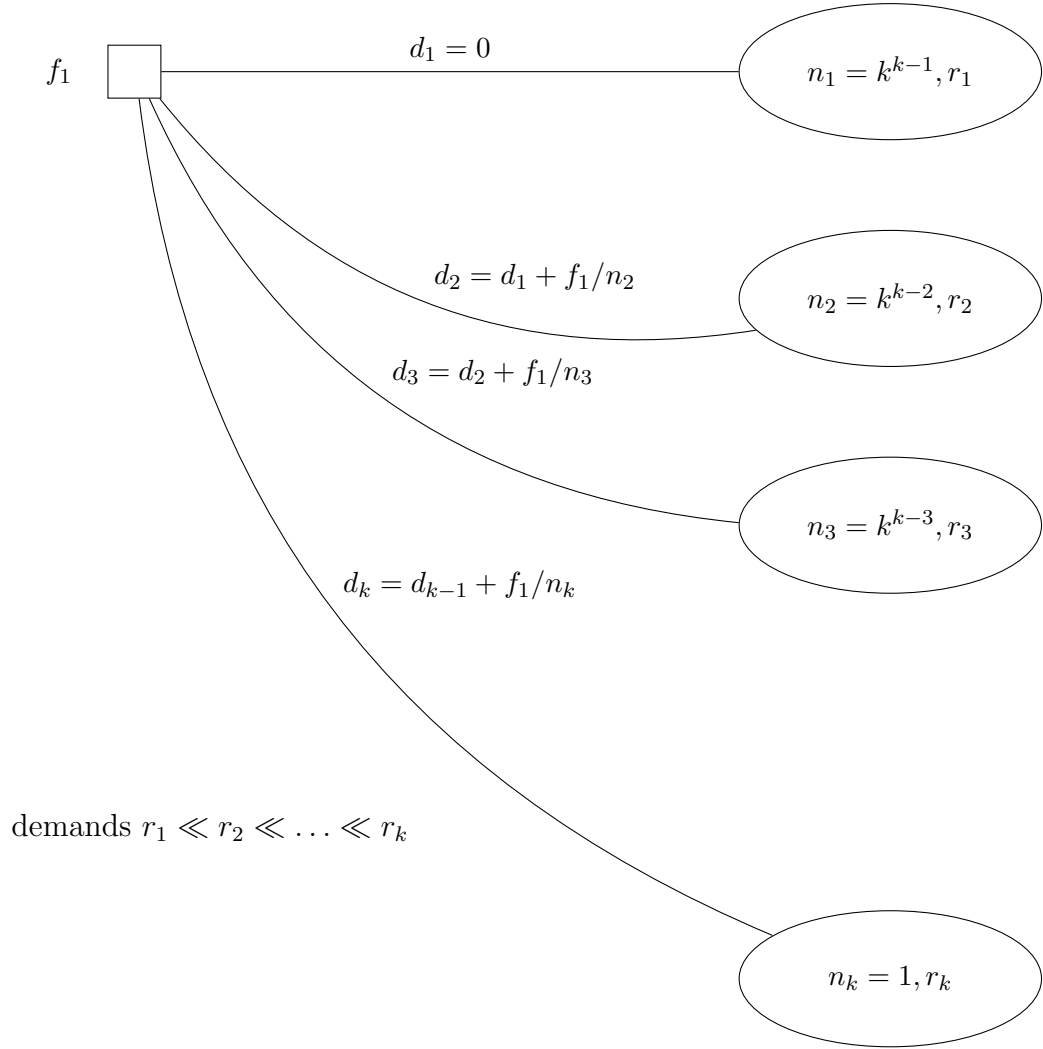


Figure 6.1: An example showing the greedy algorithm for FTFP, analyzed using dual-fitting, could give a solution with cost $\Omega(\log n / \log \log n)$ from the optimal value, assuming facility cost can only be charged to clients within the star.

For the numbers, we need $r_1 \ll r_2 \ll \dots \ll r_k$, and $n_1 = u^{k-1}, n_2 = u^{k-2}, \dots, n_k = u^0 = 1$ for some number u (Actually we take $u = k$, this choice may not be the best possible).

Call a star with facility cost zero *trivial*. It is *non-trivial* if the facility has non-zero cost. Now the greedy execution goes like this: the first non-trivial star (with r_1 replica) is (f_1, n_1) . Then we have a trivial star of zero cost facility and all n_2 clients in group 2 for r_1 replica. The second non-trivial star (with r_2 replica) is (f_1, n_2) . Notice that $r_2 \gg r_1$. The r_1 replica of trivial star with group 2 satisfy r_1 demand of the n_2 group. After that the n_2 group clients each has residual demand $r_2 - r_1 = r_2$. The process repeats until the k^{th} group finishes with r_k new facilities.

According to our local charging assumption, we have $\alpha_1 = f_1$, now defined as the total dual value of clients in group n_1 , regardless how the analysis would distribute within that group. Similarly $\alpha_2 = f_1 + n_2 d_2$, and so on. Substituting in the numbers, we have

$$\alpha_1 = f_1$$

$$\alpha_2 = f_1 + n_2 d_2 = f_1 + f_1/n_1 \cdot n_2 = f_1(1 + n_2/n_1)$$

$$\alpha_3 = f_1 + n_3 d_3 = f_1 + f_1\left(\frac{1}{n_2} + \frac{1}{n_1}\right)n_3 = f_1\left(1 + \frac{n_3}{n_2} + \frac{n_3}{n_1}\right)$$

$$\dots$$

$$\alpha_k = f_1 + n_k d_k = f_1 + f_1 n_k \left(\frac{1}{n_{k-1}} + \dots + \frac{1}{n_1}\right)$$

Notice that $r_1 \ll r_2 \ll \dots \ll r_k$ implies α_j is decided by the max among α_j^l .

Now back to the dual constraint, it requires that the shrinking factor γ needs to

satisfy the following inequality:

$$\frac{\alpha_1}{\gamma} - d_1 + \frac{\alpha_2}{\gamma} - d_2 + \dots + \frac{\alpha_k}{\gamma} - d_k \leq f_1. \quad (6.5)$$

Substitute in the α_j values derived above, we have

$$\begin{aligned} \gamma &\geq (\sum_{j=1}^k \alpha_j) / (f_1 + \sum_{j=1}^k d_j) \\ &\geq \frac{f_1 + n_1 d_1 + f_1 + n_2 d_2 + f_1 + n_3 d_3 + \dots + f_1 + n_k d_k}{f_1 + n_1 d_1 + n_2 d_2 + \dots + n_k d_k} \\ &= 1 + (k-1) f_1 / (f_1 + n_1 d_1 + n_2 d_2 + \dots + n_k d_k) \\ &= 1 + (k-1) f_1 / \left(f_1 + n_2 f_1 / n_1 + \dots + n_k f_1 \left(\frac{1}{n_{k-1}} + \frac{1}{n_{k-2}} + \dots + \frac{1}{n_1} \right) \right) \\ &= 1 + (k-1) / \left(1 + n_2 / n_1 + \dots + n_k \left(\frac{1}{n_{k-1}} + \frac{1}{n_{k-2}} + \dots + \frac{1}{n_1} \right) \right) \\ &= 1 + (k-1) / \left(1 + 1/u + \dots + \left(\frac{1}{u} + \dots + \frac{1}{u^{k-1}} \right) \right) \\ &= 1 + (k-1) / \left(1 + k/u + (k-1)/u^2 + \dots + 1/u^{k-1} \right) \\ &\geq 1 + (k-1) / \left(1 + k/u + k/u^2 + \dots + k/u^{k-1} \right) \\ &= 1 + (k-1) / \left(1 + 1 + 1/k + \dots + 1/k^{k-2} \right) \\ &\approx k/2 \end{aligned}$$

So for k groups we can force a shrinking factor γ as big as $k/2$. Recall that we have greedy being no more than H_n -approximation. Is that a contradiction? No, because we have the number of clients $n = k^{k-1} + k^{k-2} + \dots + 1 = k^k$, so $k = O(\log n / \log \log n)$. Therefore, the example shows that dual fitting with local charging cannot hope to get $O(\log n / \log \log n)$ ratio or better.

Remark Notice this example is similar in spirit to the $\Omega(\log n / \log \log n)$ example for Hochbaum's algorithm for UFL, constructed by Mahdian *et al.* [18].

Chapter 7

Conclusion

In this dissertation we studied the fault-tolerant facility placement problem (FTFP), a generalization of the well-known uncapacitated facility location problem (UFL). We demonstrated that the known LP-rounding algorithms for UFL can be adapted to FTFP while preserving the approximation ratio. To accomplish this reduction, we developed two techniques, namely demand reduction and adaptive partition, which could be of more general interest. Our results demonstrated that FTFP seems easier to approximate, compared to FTFL.

We also studied the primal-dual and dual-fitting approach, and provided a possible explanation on the difficult to obtain a constant approximation ratio using those techniques.

We hope our work in this dissertation will help other researchers interested in the fault-tolerant variant of the facility location problems to develop more insight into the difficulty and possible solutions when clients demand more than one facility and we still need to keep total cost under control.

In anticipating future research, we tend to agree with the authors, Byrka *et al.*, with their remark on UFL and FTFL, that both problems are likely to have approximation algorithms with ratio matching the 1.463 lower bound. And from our demand reduction technique, it is almost surely that FTFP shall have a 1.463-approximation algorithm, provided that FTFL can be approximated to meet the lower bound.

Bibliography

- [1] Vijay Arya, Naveen Garg, Rohit Khandekar, Adam Meyerson, Kamesh Munagala, and Vinayaka Pandit. Local search heuristic for k-median and facility location problems. In *Proceedings of the 33rd ACM Symposium on Theory of Computing, STOC '01*, pages 21–29, 2001.
- [2] M. L. Balinski. On finding integer solutions to linear programs. In *Proceedings of the IBM Scientific Computing Symposium on Combinatorial Problems*, pages 225–248, 1966.
- [3] Adriana Bumb. *Approximation Algorithms for Facility Location Problems*. PhD thesis, University of Twente, the Netherlands, 2002.
- [4] Jaroslaw Byrka. An optimal bifactor approximation algorithm for the metric uncapacitated facility location problem. In *Proceedings of the 10th International Workshop on Approximation and the 11th International Workshop on Randomization, and Combinatorial Optimization. Algorithms and Techniques, APPROX '07/RANDOM '07*, pages 29–43, 2007.
- [5] 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.
- [6] Jaroslaw Byrka, MohammadReza Ghodsi, and Aravind Srinivasan. LP-rounding algorithms for facility-location problems. *CoRR*, abs/1007.3611, 2010.
- [7] 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.
- [8] Moses Charikar and Sudipto Guha. Improved combinatorial algorithms for facility location problems. *SIAM J. Comput.*, 34(4):803–824, 2005.
- [9] Fabián Chudak and David Shmoys. Improved approximation algorithms for the uncapacitated facility location problem. *SIAM J. Comput.*, 33(1):1–25, 2004.

- [10] Fabián A. Chudak. Improved approximation algorithms for uncapacitated facility location. In *Proceedings of the 6th Integer Programming and Combinatorial Optimization, IPCO '98*, pages 180–194, 1998.
- [11] Vašek Chvátal. *Linear Programming*. W. H. Freeman and Company, 1983.
- [12] Uriel Feige. A threshold of $\ln n$ for approximating set cover. *Journal of the ACM*, 45(4):634–652, 1998.
- [13] 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.
- [14] Sudipto Guha, Adam Meyerson, and Kamesh Munagala. A constant factor approximation algorithm for the fault-tolerant facility location problem. *J. Algorithms*, 48(2):429–440, 2003.
- [15] Anupam Gupta. Lecture notes: CMU 15-854b, spring 2008, 2008.
- [16] Godfrey Harold Hardy, John Edensor Littlewood, and George Pólya. *Inequalities*. Cambridge University Press, 1988.
- [17] Dorit S. Hochbaum. Heuristics for the fixed cost median problem. *Mathematical Programming*, 22:148–162, 1982.
- [18] 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.
- [19] Kamal Jain, Mohammad Mahdian, and Amin Saberi. A new greedy approach for facility location problems. In *Proceedings of the 34th ACM Symposium on Theory of Computing, STOC '02*, pages 731–740, 2002.
- [20] 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.
- [21] Kamal Jain and Vijay V. Vazirani. An approximation algorithm for the fault tolerant metric facility location problem. *Algorithmica*, 38(3):433–439, 2003.
- [22] 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.
- [23] Mohammad Mahdian, Evangelos Markakis, Amin Saberi, and Vijay Vazirani. A greedy facility location algorithm analyzed using dual fitting. In *Proc. 4th International Workshop on Approximation Algorithms for Combinatorial Optimization Problems (APPROX)*, pages 127–137, London, UK, 2001. Springer-Verlag.

- [24] Richard L. Francis Pitu B. Mirchandani. *Discrete Location Theory*. Wiley-Interscience, 1990.
- [25] 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.
- [26] 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.
- [27] Chaitanya Swamy and David Shmoys. Fault-tolerant facility location. *ACM Trans. Algorithms*, 4(4):1–27, 2008.
- [28] Jens Vygen. Approximation algorithms for facility location problems (lecture notes. Technical Report No. 05950, Research Institute for Discrete Mathematics, University of Bonn, 2005.
- [29] 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.
- [30] Li Yan and Marek Chrobak. Approximation algorithms for the fault-tolerant facility placement problem. *Inf. Process. Lett.*, 111(11):545–549, 2011.
- [31] Li Yan and Marek Chrobak. Lp-rounding algorithms for the fault-tolerant facility placement problem. In *Proceedings of the 8th International Conference on Algorithms and Complexity, CIAC '13*, pages 370–381, 2013.

Appendix A

Technical Background

A.1 Linear Programming and Integer Programming

In this section we give a short introduction on Linear Programming and Integer Programming with an emphasis on their applicability in design and analysis of approximation algorithms for NP-optimization problems.

Most NP-optimization problems have a natural integer program in which we use variables to represent parameters in the solution that we seek, and write the constraints imposed by feasibility of the solution. The objective function is obtained by the cost function of the solution, specified by the problem. For example, in the Vertex Cover problem, we are given a graph $G = (V, E)$ and we are to find a subset W of V , such that every edge $e \in E$ has at least one endpoint in W , and we want the set W to have minimum size. To formulate this problem as an integer program, we use $x_v \in \{0, 1\}$ to denote whether a node $v \in V$ is in W or not. The constraint is that for every edge $e = (u, v)$, we have $x_u + x_v \geq 1$.

The objective is to minimize $\sum_{v \in V} x_v$. The integer program for Vertex Cover is written as

$$\begin{aligned} & \text{minimize } \sum_{v \in V} x_v \\ & \text{subject to } x_u + x_v \geq 1 \quad \forall (u, v) \in E \\ & x_v \in \{0, 1\} \quad \forall v \in V \end{aligned}$$

In general an integer program cannot be solved exactly in polynomial time, as integer programming is **NP**-hard. However, if we relax the integral constraint and allow the variables to take fractional value, we then obtain a Linear Program (LP) and LP is polynomially solvable, for example, using the ellipsoid method or the interior point method. Thus we can first solve the LP optimally, obtaining a fractional optimal solution to the LP. The value of the fractional optimal solution is then a lower bound on the value of the optimal integral solution, assuming a minimization problem. Our next step is then to round the fractional solution appropriately, so that we maintain the feasibility while keeping the cost from increasing too much. The exact rounding procedure is problem specific and we shall not delve into the details here. The rounding relevant to the FTFP problem in this thesis is presented in detail in Chapter 5.

We now give a brief introduction of linear programming, see [11] for an introductory

book on this topic. A general Linear Program can be written as

$$\begin{aligned}
& \text{minimize} && \sum_{j=1}^n c_j x_j && (A.1) \\
& \text{subject to} && \sum_{j=1}^n a_{ij} x_j \geq b_i, && \text{for } i = 1, \dots, m \\
& && x_j \geq 0 && \text{for } j = 1, \dots, n
\end{aligned}$$

For the LP above, we can take its dual as

$$\begin{aligned}
& \text{maximize} && \sum_{i=1}^m b_i y_i && (A.2) \\
& \text{subject to} && \sum_{i=1}^m a_{ij} y_i \leq c_j && \text{for } j = 1, \dots, n \\
& && y_i \geq 0 && \text{for } i = 1, \dots, m
\end{aligned}$$

The LP (A.1) is called the primal program and the LP (A.2) is called the dual program.

The weak duality theorem tells us that, for every feasible solution \mathbf{x} for the primal (A.1)

and \mathbf{y} for the dual (A.2), we have that $\mathbf{c}^T \mathbf{x} \geq \mathbf{b}^T \mathbf{y}$. The strong duality theorem tells

us that if both the primal (A.1) and the dual (A.2) are feasible, then both of them have

optimal solution \mathbf{x}^* and \mathbf{y}^* and their objective function values equal, that is $\mathbf{c}^T \mathbf{x}^* = \mathbf{b}^T \mathbf{y}^*$.

Moreover, the complementary slackness conditions assert that two feasible solutions \mathbf{x} and \mathbf{y}

are both optimal to LP (A.1) and (A.2) respectively, if and only if, for every primal variable

x_j , either $x_j = 0$ or the corresponding constraint in the dual is tight, that is $\sum_{i=1}^m a_{ij} y_i = c_j$.

And for every dual variable y_i , either $y_i = 0$ or the corresponding constraint in the primal

is tight, that is $\sum_{j=1}^n a_{ij}x_j = b_i$. The complementary slackness conditions provide a simple way to validate the optimality when one is presented with two solutions, proposed to be optimal for the primal and dual program respectively.

The complementary slackness conditions play a crucial role in the design and analysis of approximation algorithms. For example, suppose we have an algorithm that computes a feasible integral solution \mathbf{x} to the primal program (A.1) and a feasible integral solution to the dual program (A.2). Moreover, we know that the two solutions satisfy a relaxed version of the complementary slackness conditions: for some numbers γ and ρ , we have

$$\begin{aligned} \text{either } y_i = 0 \quad \text{or} \quad b_i &\leq \sum_j a_{ij}x_j \leq \gamma b_i & \text{for } i = 1, \dots, m. \\ \text{either } x_j = 0 \quad \text{or} \quad \rho c_j &\leq \sum_i a_{ij}y_i \leq c_j & \text{for } j = 1, \dots, n. \end{aligned}$$

Then the integral solution \mathbf{x} has cost no more than γ/ρ times the optimal value. In particular, we have $\sum_j c_j x_j \leq \gamma/\rho \sum_i b_i y_i$ and the value for a feasible dual solution, namely $\sum_i b_i y_i$, is a lower bound on the optimal value of the primal program.

As an application of the complementary slackness conditions, we look at their use in the design and analysis of algorithms for the facility location problems. Recall that we define the neighborhood $N(j)$ of a client j as the set of facilities with $x_{ij}^* > 0$, where $\mathbf{x}^*, \mathbf{y}^*$ is some fractional optimal fractional solution and $\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*$ is some optimal fractional dual solution. The complementary slackness conditions provide an upper bound on the maximum distance from a facility $i \in N(j)$ to a client j , since one dual constraint says

$\alpha_j - \beta_{ij} \leq d_{ij}$ and if the primal solution has $x_{ij}^* > 0$, then the inequality is actually an equality and we have $\alpha_j^* - \beta_{ij}^* = d_{ij}$. Together with $\beta_{ij}^* \geq 0$, we have $\alpha_j^* \geq d_{ij}$ for every i such that $x_{ij}^* > 0$.

The idea of using relaxed complementary slackness conditions in designing approximation algorithms for the uncapacitated facility location problem is demonstrated by Jain and Vazirani [20]. They proposed an algorithm that outputs an integral solution (\mathbf{x}, \mathbf{y}) to the primal program (3.1) and a feasible (possibly fractional) solution $(\boldsymbol{\alpha}, \boldsymbol{\beta})$ to the dual program (3.2)¹. Moreover, the two solutions satisfy the conditions that

$$\text{either } \sum_j \beta_{ij} = f_i \quad \text{or} \quad y_i = 0.$$

$$\text{either } 1/3 d_{ij} \leq \alpha_j - \beta_{ij} \leq d_{ij} \quad \text{or} \quad x_{ij} = 0.$$

The solution (\mathbf{x}, \mathbf{y}) then is an 3-approximation to the optimal solution.

A.2 Proof of Inequality (5.3)

In Sections 5.2 and 5.3 we use the following inequality

$$\begin{aligned} \bar{d}_1 g_1 + \bar{d}_2 g_2 (1 - g_1) + \dots + \bar{d}_k g_k (1 - g_1)(1 - g_2) \dots (1 - g_{k-1}) \\ \leq \frac{1}{\sum_{s=1}^k g_s} \left(\sum_{s=1}^k \bar{d}_s g_s \right) \left(\sum_{t=1}^k g_t \prod_{z=1}^{t-1} (1 - g_z) \right). \end{aligned} \tag{A.3}$$

for $0 < \bar{d}_1 \leq \bar{d}_2 \leq \dots \leq \bar{d}_k$, and $0 < g_1, \dots, g_s \leq 1$.

We give here a new proof of this inequality, much simpler than the existing proof

¹For the uncapacitated facility location problem we have all $r_j = 1$ for $j \in \mathbb{C}$.

in [9], and also simpler than the argument by Sviridenko [26]. We derive this inequality from the following generalized version of the Chebyshev Sum Inequality:

$$\sum_i p_i \sum_j p_j a_j b_j \leq \sum_i p_i a_i \sum_j p_j b_j, \quad (\text{A.4})$$

where each summation 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 (A.4), we can obtain our inequality (A.3) by simple substitution

$$p_i \leftarrow g_i, a_i \leftarrow \bar{d}_i, b_i \leftarrow \prod_{s=1}^{i-1} (1 - g_s),$$

for $i = 1, \dots, k$.

For the sake of completeness, we include the proof of inequality (A.4), due to Hardy, Littlewood and Polya [16]. 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.