# Efficient Solution of a Class of Location-Allocation Problems with Stochastic Demand and Congestion

Navneet Vidyarthi[a,*], Sachin Jayaswal[b]

[a]Department of Supply Chain and Business Technology Management,
John Molson School of Business, Concordia University, Montreal, QC, H3G 1M8, Canada

[b]Production and Quantitative Methods, Indian Institute of Management,
Vastrapur, Ahmedabad, Gujarat, 380 015, India

## Abstract

We consider a class of location-allocation problems with immobile servers, stochastic demand and congestion that arises in several planning contexts: location of emergency medical clinics; preventive healthcare centers; refuse collection and disposal centers; stores and service centers; bank branches and automated teller machines; internet mirror sites; and distribution centers in supply chains. The problem seeks to simultaneously locate service facilities, equip them with appropriate capacities, and allocate customer demand to these facilities such that the total cost, which consists of the fixed cost of opening facilities with sufficient capacities, the access cost of users' travel to facilities, and the queuing delay cost, is minimized. Under Poisson user demand arrivals and general service time distributions, the problem is set up as a network of independent $M/G/1$ queues, whose locations, capacities and service zones need to be determined. The resulting mathematical model is a non-linear integer program. Using simple transformation and piecewise linear approximation, the model is linearized and solved to $\epsilon$-optimality using a constraint generation method. Computational results are presented for instances up to 400 users, 25 potential service facilities, and 5 capacity levels with different coefficient of variation of service times and average queueing delay costs per customer. The results indicate that the proposed solution method is efficient in solving a wide range of problem instances.

*Keywords:* Service System Design; Location-Allocation; Queueing; Stochastic Demand; Congestion; Constraint Generation Method

## 1. Introduction

Problems arising in several planning contexts require deciding: (i) the location of service facilities and their capacities; and (ii) service zones (allocations) of the located service facilities. Examples include location-allocation of emergency service facilities such as medical clinics and preventive health care facilities (Zhang et al., 2009, 2010, 2012); stores and service centers; bank branches and automated teller machines (Aboolian et al., 2008; Wang et al., 2002); internet mirror sites; and distribution centers in supply chains (Huang et al., 2005; Vidyarthi et al., 2009). All the above examples are characterized by servers (medical clinics, bank branches, distribution centers, etc.) that are immobile in that the customers

---

*Corresponding author, Phone: +001-514-848-2424x2990, Fax: +001-514-848-2824
  *Email addresses:* `navneetv@jmsb.concordia.ca` (Navneet Vidyarthi), `sachin@iimahd.ernet.in` (Sachin Jayaswal)

need to travel to the service facilities to avail of their services, as opposed to the servers travelling (mobile servers) to the customers' site in response to calls for their services. Such problems are generally also characterized by random (stochastic) nature of service calls (demand arrivals) and their service requirements (service times). These problems are commonly known in the literature as facility location problems with immobile servers, stochastic demand and congestion (Berman and Krass, 2004). They are also termed as service system design problems with stochastic demand and congestion (Amiri, 1997, 1998, 2001; Elhedhli, 2006). Excellent reviews on this class of problems are provided by Berman and Krass (2004) and Boffey et al. (2007).

For facility location problems with stochastic demand and congestion, the following two factors are important: (i) the costs of providing service; and (ii) the quality of service, with an objective generally requiring a balance between the two. The costs of providing service are related to the fixed cost of opening/operating the service facilities and the cost of accessing these facilities by the users. The service quality, on the other hand, is often measured in terms of: (i) the average number of users waiting for service; (ii) average waiting time per user; or (iii) the probability of serving a user within a time limit (Elhedhli, 2006). Balance between service costs and service quality is commonly achieved in the literature using a combination of the total cost of opening and accessing facilities and the cost associated with waiting customers, which is minimized in the objective function (Amiri, 1997, 1998; Wang et al., 2002; Elhedhli, 2006; Castillo et al., 2009). Others in the literature minimize the cost of providing service subject to a minimum threshold on the service quality, where the service quality may be defined in one of the ways described above (Marianov and Serra, 1998, 2002; Silva and Serra, 2008).

In the current work, we use the former of the two approaches described above, i.e., we consider as an objective the minimization of a combination of the total cost of opening and accessing facilities and the cost associated with waiting customers. We note that due to the complexity of the underlying problem, most papers in this category make assumptions such as: (i) either the number or capacity of the facilities (or both) are fixed; (ii) the assignment of users to the facilities are known in advance (closest assignment property); (iii) the demand arrival process is Poisson; and (iv) the service times follow an exponential distribution (see Amiri, 1997; Marianov and Serra, 2002; Wang et al., 2002; Elhedhli, 2006; Aboolian et al., 2008, and references therein). Despite these simplifying assumptions, the techniques proposed to date to solve the problem, with the exception of Elhedhli (2006), are either approximate or heuristic based.

The contribution of this paper is two fold. First, we present a more generalized model of the problem than the extant literature by assuming a general distribution for the service times at facilities, as opposed to exponential distribution used in the literature. More specif-

ically, our proposed model seeks to determine the minimum cost configuration (location of service facilities with adequate capacity and allocation of service zones to these facilities) of a service system under Poisson arrivals and general service time distribution, where the total cost consists of the costs of opening and accessing service facilities and the cost associated with waiting customers. The proposed model, therefore, is more challenging to solve than the ones available in the literature that assume exponential service time distribution, which themselves are too difficult to solve using exact methods. So, our second contribution lies in the exact ($\epsilon$-optimal) solution method that we propose to solve our model. Our proposed solution method is based on a simple transformation and piecewise linearization of our non-linear integer programming (IP) model, which is solved to optimality (or $\epsilon$-optimality) using a constraint generation algorithm.

The remainder of the paper is organized as follows. In Section 2, we describe the problem setting, followed by its non-linear IP model. Section 3 describes the transformation and the piecewise linearization approach for the non-linear IP model. To solve the linearized model, we present a constraint generation based solution approach in Section 4. Computational results are reported in Section 5. Section 6 concludes with some directions for future research.

## 2. Problem Formulation

Consider a set of user nodes, each indexed by $i \in I$ whose demand for service occurs continuously over time according to an independent Poisson process with rate $\lambda_i$. We consider a directed choice environment, where users are assigned to facilities, each indexed by $j \in J$, by a central decision maker. We assume that users from any node are entirely assigned to a single service facility, where each facility operates as a single server with an infinite buffer to accommodate users waiting for service. If $x_{ij}$ is a binary variable that equals 1 if the demand for service from user node $i$ is satisfied by facility $j$, and 0 otherwise, then the aggregate demand arrival rate at facility $j$, as a result of the superposition of Poisson processes, also follows a Poisson process with mean $\Lambda_j = \sum_{i \in I} \lambda_i x_{ij}$ (Gross and Harris, 1998).

Let $y_{jk}$ be a binary variable that equals 1 if facility at site $j$ is open and equipped with a capacity level $k \in K$, 0 otherwise. Further, assume that the service times at any facility $j$ are independent and identically distributed with a mean $1/\mu_{jk}$ and variance $\sigma_{jk}^2$ if it is equipped with a capacity level $k$. Any facility $j$ is thus modeled as an $M/G/1$ queue with a service rate $\mu_j = \sum_{k \in K} \mu_{jk} y_{jk}$ and variance in service times given by $\sigma_j^2 = \sum_{k \in K} \sigma_{jk}^2 y_{jk}$. Thus, the service system design problem is modeled as a network of independent M/G/1 queues.

Under steady state conditions ($\Lambda_j/\mu_j < 1$), first-come-first-serve (FCFS) queuing discipline, and infinite buffers to accommodate users waiting for service, the expected waiting

time (including the time spent in service) of users at facility $j$ is given, by the Pollaczek-Khintchine formula, (Gross and Harris, 1998) as:

$$E[w_j] = \left(\frac{1 + Cv_j^2}{2}\right)\frac{\tau_j \rho_j}{1 - \rho_j} + \tau_j = \left(\frac{1 + Cv_j^2}{2}\right)\frac{\Lambda_j}{\mu_j(\mu_j - \Lambda_j)} + \frac{1}{\mu_j} \qquad (1)$$

where $\tau_j = 1/\mu_j$ is the average service time at facility $j$, $\rho_j = \Lambda_j/\mu_j$ is the average utilization of facility $j$, and $Cv_j = \sigma_j \mu_j$ is the coefficient of variation of service times at facility $j$. $E[w_j]$ can be written in terms of location and allocation variables ($y_{jk}$ and $x_{ij}$) as:

$$E[w_j(\mathbf{x}, \mathbf{y})] = \frac{\left(1 + \sum_{k \in K} Cv_{jk}^2 y_{jk}\right) \sum_{i \in I} \lambda_i x_{ij}}{2 \sum_{k \in K} \mu_{jk} y_{jk} \left(\sum_{k \in K} \mu_{jk} y_{jk} - \sum_{i \in I} \lambda_i x_{ij}\right)} + \frac{1}{\sum_{k \in K} \mu_{jk} y_{jk}} \qquad (2)$$

The expected number of users in service or waiting for service at facility $j$ is given, using Little's law, as $\Lambda_j E[w_j]$. If $d$ denotes the average waiting time cost per customer (henceforth called unit queuing delay cost), then the *total delay/congestion cost* in the network can be expressed as $d \sum_{j \in J} \Lambda_j E[w_j(\mathbf{x}, \mathbf{y})] = \sum_{j \in J} \sum_{i \in I} \lambda_i x_{ij} E[w_j(\mathbf{x}, \mathbf{y})]$. We assume there is a fixed set up cost $f_{jk}$ (amortized over the planning period) of locating a facility with capacity level $k$ at site $j$, and a variable access cost $c_{ij}$ of providing service to users at node $i$ from facility at site $j$. The problem is to simultaneously determine: (i) the locations of the service facilities and their corresponding capacity levels; (ii) the assignment of users to located service facilities, such that the total system-wide cost, consisting of cost of opening service facilities with appropriate capacities, cost of accessing service facilities by users and cost associated with customers' waiting, is minimized. The resulting non-linear integer program (IP) model of the problem is as follows:

$$[P]: Z(\mathbf{x}, \mathbf{y}) = \min_{\mathbf{x}, \mathbf{y}} \sum_{j \in J} \sum_{k \in K} f_{jk} y_{jk} + \sum_{i \in I} \sum_{j \in J} c_{ij} x_{ij} + d \sum_{j \in J} \sum_{i \in I} \lambda_i x_{ij} E[w_j(\mathbf{x}, \mathbf{y})] \qquad (3)$$

$$\text{s.t.} \quad \sum_{i \in I} \lambda_i x_{ij} \leq \sum_{k \in K} \mu_{jk} y_{jk} \qquad \forall j \qquad (4)$$

$$\sum_{j \in J} x_{ij} = 1 \qquad \forall i \qquad (5)$$

$$\sum_{k \in K} y_{jk} \leq 1 \qquad \forall j \qquad (6)$$

$$x_{ij}, y_{jk} \in \{0, 1\} \qquad \forall i, j, k \qquad (7)$$

The three terms in the objective function (3) are: (i) cost of opening facilities with appropriate capacities; (ii) cost of accessing service facilities; and (iii) cost of users' waiting at facilities. The expression for $E[w_j(\mathbf{x}, \mathbf{y})]$ in the objective function is given by (2). Constraint set (4) ensures that at every facility, the total demand allocated is less than its capacity. Note that constraint set (4) will be non-binding at optimality, else the term $E[w_j(\mathbf{x}, \mathbf{y})]$ in the objective function goes to infinity. This ensures the stability of the queueing system ($\rho_j = \Lambda_j/\mu_j < 1$) at each facility $j$. Constraint set (5) ensures that each user node is

4

assigned to only one of the open facilities for its service. Constraint set (6) states that at most one among multiple capacity levels is selected at a facility. Constraint set (7) imposes binary restrictions on the location and allocation variables.

The presence of the non-linear term $\sum_{i \in I} \sum_{j \in J} \lambda_i x_{ij} E[w_j(\mathbf{x}, \mathbf{y})]$ in the objective function makes $[P]$ challenging to solve. In the next section, we present an approach to linearize the expression for the total waiting time spent by the users at a facility, followed by an exact solution procedure, based on a constraint generation algorithm, to solve the linearized model.

## 3. Model Linearization

The non-linear term in the objective function of $[P]$ can be written, using (1), as:

$$\Lambda_j E[w_j] = \left( \frac{1 + Cv_j^2}{2} \right) \frac{\Lambda_j^2}{\mu_j(\mu_j - \Lambda_j)} + \frac{\Lambda_j}{\mu_j} = \left( \frac{1 + Cv_j^2}{2} \right) \frac{\rho_j^2}{(1 - \rho_j)} + \rho_j \tag{8}$$

To linearize (8), it can be rewritten, upon rearranging its terms, as:

$$\Lambda_j E[w_j] = \frac{1}{2} \left\{ \left( 1 + Cv_j^2 \right) \frac{\rho_j}{1 - \rho_j} + \left( 1 - Cv_j^2 \right) \rho_j \right\} \tag{9}$$

Let us define a set of nonnegative auxiliary variables, $U_j$, such that:

$$U_j = \frac{\rho_j}{1 - \rho_j} = \frac{\Lambda_j}{\mu_j - \Lambda_j} = \frac{\sum_{i \in I} \lambda_i x_{ij}}{\sum_{k \in K} \mu_{jk} y_{jk} - \sum_{i \in I} \lambda_i x_{ij}} \tag{10}$$

which implies:

$$\rho_j = \frac{U_j}{1 + U_j} \tag{11}$$

Using $\rho_j = \Lambda_j / \mu_j$, the total demand $\Lambda_j$ at facility $j$ can be expressed as:

$$\Lambda_j = \sum_{i \in I} \lambda_i x_{ij} = \rho_j \mu_j = \rho_j \sum_{k \in K} \mu_{jk} y_{jk} = \sum_{k \in K} \mu_{jk} z_{jk}, \quad \text{where } z_{jk} = \begin{cases} \rho_j, & \text{if } y_{jk} = 1 \\ 0, & \text{otherwise} \end{cases}$$

Hence, the total demand at any facility $j$ can be expresses as:

$$\sum_{i \in I} \lambda_i x_{ij} = \sum_{k \in K} \mu_{jk} z_{jk} \qquad \forall j$$

We know that any facility can have at most one capacity level, i.e., there exists at most one $k = k'$ such that $y_{jk'} = 1$, while $y_{jk} = 0 \ \forall \ k \neq k'$. Further, $\rho_j < 1$. Using this knowledge,

$z_{jk}$ can alternatively be expressed using the following set of constraints:

$$z_{jk} \leq y_{jk} \qquad\qquad \forall j, k$$

$$\sum_{k \in K} z_{jk} = \rho_j \qquad\qquad \forall j$$

$$z_{jk} \geq 0 \qquad\qquad \forall j, k$$

With the above substitutions in (9), the expression for $\Lambda_j E[w_j]$ reduces to:

$$
\begin{aligned}
\Lambda_j E[w_j] &= \frac{1}{2}\left\{ \left(1 + \sum_{k \in K} Cv_{jk}^2 y_{jk}\right) U_j + \left(1 - \sum_{k \in K} Cv_{jk}^2 y_{jk}\right) \rho_j \right\} \\
&= \frac{1}{2}\left( U_j + \sum_{k \in K} Cv_{jk}^2 w_{jk} + \rho_j - \sum_{k \in K} Cv_{jk}^2 z_{jk} \right), \qquad \text{where } w_{jk} = \begin{cases} U_j, & \text{if } y_{jk} = 1 \\ 0, & \text{otherwise} \end{cases}
\end{aligned}
$$

Again, using the fact that there exists at most one $k = k'$ such that $y_{jk'} = 1$, while $y_{jk} = 0 \ \forall \ k \neq k'$, $w_{jk}$ can alternatively be expressed using the following set of constraints:

$$w_{jk} \leq My_{jk} \qquad\qquad \forall j, k$$

$$\sum_{k \in K} w_{jk} = U_j \qquad\qquad \forall j$$

$$w_{jk} \geq 0 \qquad\qquad \forall j, k$$

where $M$ is a large number (Big-M). We now state a Lemma that helps us linearize the above non-linear model $[P]$.

**Lemma 1:** *The function $\rho_j(U_j) = \frac{U_j}{1+U_j}$ is concave in $U_j \in [0, \infty)$.*
**Proof:**
Differentiating $\rho_j$ w.r.t. $U_j$, we get the first derivative, $\frac{\delta \rho_j}{\delta U_j} = \frac{1}{(1+U_j)^2} > 0$, and the second derivative, $\frac{\delta^2 \rho_j}{\delta U_j^2} = \frac{-2}{(1+U_j)^3} < 0$, which proves that the function $\rho_j(U_j)$ is concave in $U_j$. ∎

Lemma 1 implies that for a given set of points indexed by $h$, $h \in H$, the function $\rho_j(U_j)$ can be approximated arbitrarily close by a set of piecewise linear functions that are tangent to $\rho_j$ at points $\{U_j^h\}_{h \in H}$, such that:

$$\rho_j = \min_{h \in H}\left\{ \frac{1}{(1+U_j^h)^2} U_j + \frac{(U_j^h)^2}{(1+U_j^h)^2} \right\} \qquad\qquad (12)$$

This is equivalent to the following set of constraints:

$$\rho_j \leq \frac{1}{(1+U_j^h)^2} U_j + \frac{(U_j^h)^2}{(1+U_j^h)^2}, \qquad\qquad \forall j, h \in H \qquad (13)$$

provided $\exists\, h \in H$ such that (13) holds with equality.

Using the above substitutions result in the following linear mixed integer program (MIP) reformulation of [P]:

$$[P(H)] : \min \sum_{j \in J} \sum_{k \in K} f_{jk} y_{jk} + \sum_{i \in I} \sum_{j \in J} c_{ij} x_{ij} + \frac{d}{2} \sum_{j \in J} \left\{ U_j + \rho_j + \sum_{k \in K} Cv_{jk}^2 (w_{jk} - z_{jk}) \right\} \quad (14)$$

$$\text{s.t. } (5) - (7), (13)$$

$$\sum_{i \in I} \lambda_i x_{ij} - \sum_{k \in K} \mu_{jk} z_{jk} = 0 \qquad \forall j \tag{15}$$

$$z_{jk} \leq y_{jk} \qquad \forall j,k \tag{16}$$

$$\sum_{k \in K} z_{jk} = \rho_j \qquad \forall j \tag{17}$$

$$w_{jk} \leq M y_{jk} \qquad \forall j,k \tag{18}$$

$$\sum_{k \in K} w_{jk} = U_j \qquad \forall j \tag{19}$$

$$0 \leq z_{jk}, \rho_j \leq 1; \quad w_{jk}, U_j \geq 0 \quad \forall j,k \tag{20}$$

Equivalence between $[P]$ and $[P(H)]$ requires that $\exists\, h \in H$ such that (13) holds with equality. Proposition 1 states that this condition will always be satisfied at optimality.

**Proposition 1:** *In the linearized model $[P(H)]$, at least one of the constraints in (13) will be binding at optimality.*

**Proof:**

Upon rearranging the terms, (13) can be rewritten as:

$$U_j \geq (1 + U_j^h)^2 \rho_j - (U_j^h)^2, \qquad \forall j, h \in H \tag{21}$$

Since $U_j$ appears in the objective function of $[P(H)]$ with a positive coefficient, $[P(H)]$ attains its minimum value only when $U_j$ is minimized. This implies that $\forall j \in J$, $\exists\, h \in H$ such that (21) holds with equality if $(1+U_j^h)^2 \rho_j - (U_j^h)^2 \geq 0$, else $U_j = 0$ if $(1+U_j^h)^2 \rho_j - (U_j^h)^2 < 0$.

Further,

$$0 \leq (1 + U_j^h)^2 \rho_j - (U_j^h)^2$$
$$= (\rho_j - 1)(U_j^h)^2 + 2\rho_j U_j^h + \rho_j$$
$$\Leftrightarrow U_j^h \in \left[ 0, \frac{\rho_j + \sqrt{\rho_j}}{1 - \rho_j} \right] \ \forall\, j \in J, h \in H \ (\text{since } \rho_j \leq 1, U_j \geq 0)$$

Thus, to prove that $\forall j \in J$, $\exists\, h \in H$ such that (21) holds with equality, it is sufficient to

show that $U_j^h \in \left[0, \frac{\rho_j + \sqrt{\rho_j}}{1 - \rho_j}\right]$. Since $U_j^h$ is an approximation to $U_j$, we obtain:

$$0 \leq U_j^h \approx U_j = \frac{\rho_j}{1 - \rho_j}$$
$$\leq \frac{\rho_j + \sqrt{\rho_j}}{1 - \rho_j}$$

This proves that $\forall j \in J, \exists h \in H$ such that (21) holds with equality. ∎

The linear MIP model $[P(H)]$ has $2(|J| + |J| * |K|)$ additional continuous variables compared to the non-linear IP model $[P]$. Further, $[P(H)]$ has $(|I| + 4*|J| + 2*|J|*|K| + |J|*|H|)$ constraints, as opposed to only $(|I| + 2*|J|)$ constraints in $[P]$. Hence, the non-linearity of $[P]$ is eliminated at the expense of having to deal with a large number of additional variables and constraints in $[P(H)]$.

### 3.1. Special Cases

In many cases, the service at facilities involve repeated steps without much variation, i.e., $Cv_{jk} = 0$ (such that each service facility is modeled as an $M/D/1$ queuing system). For such deterministic service times, the users' expected waiting time at facility $j$ is given by: $\Lambda_j E[w_j] = \frac{1}{2} \left( \frac{\Lambda_j}{\mu_j - \Lambda_j} + \frac{\Lambda_j}{\mu_j} \right) = \frac{1}{2}(U_j + \rho_j)$ and the resulting linear MIP model is as follows:

$$[P(H)_{Cv=0}] : \min \quad \sum_{j \in J} \sum_{k \in K} f_{jk} y_{jk} + \sum_{i \in I} \sum_{j \in J} c_{ij} x_{ij} + \frac{d}{2} \sum_{j \in J} (U_j + \rho_j)$$
$$\text{s.t.} \quad (4) - (7), (13)$$
$$0 \leq \rho_j \leq 1; \quad U_j \geq 0 \qquad \forall i, j, k$$

For exponentially distributed service times at the facilities, i.e., $Cv_{jk} = 1$ ($M/M/1$ case), the expression is given by: $\Lambda_j E[w_j] = \frac{\rho_j}{1 - \rho_j} = U_j$, and the linear model reduces to:

$$[P(H)_{Cv=1}] : \min \quad \sum_{j \in J} \sum_{k \in K} f_{jk} y_{jk} + \sum_{i \in I} \sum_{j \in J} c_{ij} x_{ij} + d \sum_{j \in J} U_j$$
$$\text{s.t.} \quad (4) - (7), (13)$$
$$0 \leq \rho_j \leq 1; \quad U_j \geq 0 \qquad \forall i, j, k$$

## 4. Solution Approach

We state the following two propositions, which are used in the development of the solution algorithm for $[P(H)]$.

**Proposition 2:** *For any given subset of points $\{U_j^h\}_{H^q \subset H}$, (22) provides a lower bound on the optimal objective function value of $[P]$, where $v(\bullet)$ is the objective function value of the*

8

*problem (●) and ($\mathbf{x}^p, \mathbf{y}^q, \rho^q, \mathbf{w}^q, \mathbf{z}^q, \mathbf{U}^q$) is the optimal solution to $[P(H^q)]$.*

$$LB = v(P(H^q)) =$$

$$\sum_{j \in J} \sum_{k \in K} f_{jk} y_{jk}^q + \sum_{i \in I} \sum_{j \in N} c_{ij} x_{ij}^q + \frac{d}{2} \sum_{j \in J} \left\{ U_j^q + \rho_j^q + \sum_{k \in K} Cv_{jk}^2 (w_{jk}^q - z_{jk}^q) \right\} \quad (22)$$

**Proof:**

Since $[P(H^q)]$ is a relaxation of the full problem $[P(H)]$, the objective function value of $[P(H^q)]$, given by (22), provides a lower bound on the optimal objective function value of $[P(H)]$, and hence on the optimal objective function value of $[P]$. ∎

**Proposition 3:** *For any given subset of points $\{U_j^h\}_{H^q \subset H}$, (23) provides an upper bound on the optimal objective function value of $[P]$, where ($\mathbf{x}^p, \mathbf{y}^q$) is the optimal solution to $P[(H^q)]$.*

$$UB = Z(\mathbf{x}^q, \mathbf{y}^q) = \sum_{j \in J} \sum_{k \in K} f_{jk} y_{jk}^q + \sum_{i \in I} \sum_{j \in J} c_{ij} \lambda_i x_{ij}^q +$$

$$d \left\{ \frac{\left(1 + \sum_{k \in K} Cv_{jk}^2 y_{jk}\right) \left(\sum_{i \in I} \lambda_i x_{ij}\right)^2}{2 \sum_{k \in K} \mu_{jk} y_{jk} \left(\sum_{k \in K} \mu_{jk} y_{jk} - \sum_{i \in I} \lambda_i x_{ij}\right)} + \frac{\sum_{i \in I} \lambda_i x_{ij}}{\sum_{k \in K} \mu_{jk} y_{jk}} \right\} \quad (23)$$

**Proof:**

For any subset of points $\{U_j^h\}_{H^q \subset H}$, the optimal solution ($\mathbf{x}^p, \mathbf{y}^q$) to $[P(H^q)]$ is also a feasible solution to $[P]$ as all the constraints of $[P]$ are also contained in $[P(H^q)]$. Hence, the objective function of $[P]$ evaluated at ($\mathbf{x}^p, \mathbf{y}^q$), which is given by (23), provides an upper bound on the optimal objective of $[P]$. ∎

*4.1. Solution Algorithm*

Although there are a large number of constraints/cuts (13) in the linear MIP model $[P(H)]$, it is not necessary to generate all of them. Instead, it suffices to start with a subset $H^1 \subset H$ of these cuts, where $H^1$ may be empty or chosen a priori, and generate the rest as needed. Our preliminary computational experiments, presented in Table 1, suggest a much faster convergence of the algorithm when $H^1$ is non-empty. We, therefore, use a carefully chosen subset $H^1$ of initial cuts in all our subsequent experiments. The subset of points $\{U_j^h\}_{H^1 \subset H}$, required to obtain the initial subset of cuts, is generated to approximate the function $\rho_j(U_j) = \frac{U_j}{1+U_j}$ using its tangents $\widehat{\rho}_j(U_j)$ at these points such that the approximation error, measured as $\widehat{\rho}_j(U_j) - \rho_j(U_j)$, is at most 0.001 (Elhedhli, 2005). The resulting $[P(H^1)]$ is solved, giving a solution ($\mathbf{x}^1, \mathbf{y}^1, \rho^1, \mathbf{w}^1, \mathbf{z}^1, \mathbf{U}^1$). The lower bound ($LB^1 q$) and upper bound ($UB^1$) are computed using (22) and (23) respectively. If $UB^1$ equals $LB^1$ within some accepted tolerance ($\epsilon$), then ($\mathbf{x}^1, \mathbf{y}^1$) is an optimal solution to $[P]$, and the algorithm stops. Otherwise, a new set of points $\{U_j^h\}$ is generated using the current solution as $U_j^{h new} = \frac{\sum_{i \in I} \lambda_i x_{ij}^1}{\sum_{k \in K} \mu_{jk} y_{jk}^1 - \sum_{i \in I} \lambda_i x_{ij}^1} \quad \forall j$. A new set of constraints/cuts of the form

9

(13) is generated at these new points, which are appended to $[P(H^1)]$ to give $[P(H^2)]$. Now, $[P(H^2)]$ is solved to yield a solution $(\mathbf{x}^2, \mathbf{y}^2, \rho^2, \mathbf{w}^2, \mathbf{z}^2, \mathbf{U}^2)$ and a lower bound $LB^2$. Since the upper bound, as given by (23), changes non-monotically, the new upper bound $UB^2$ is retained as $\min\{UB^1, Z(\mathbf{x}^2, \mathbf{y}^2)\}$. If $UB^2$ equals $LB^2$ within the accepted tolerance ($\epsilon$), then the algorithm terminates with $(\mathbf{x}^2, \mathbf{y}^2)$ as an optimal solution. Otherwise, the above process is repeated until $UB^q$ equals $LB^q$ within ($\epsilon$) for some iterations $q$. The details of the algorithm are outlined below:

---

**Algorithm 1** Constraint Generation Algorithm for [P(H)]

---

1: $q \leftarrow 1$; $UB^{q-1} \leftarrow +\infty$; $LB^{q-1} \leftarrow -\infty$.
2: Choose an initial set of points $\{U^h\}_{h \in H^q}$ to approximate the function $\rho_j(U_j) = U_j/1+U_j$.
3: **while** $(UB^{q-1} - LB^{q-1})/UB^{q-1} > \epsilon$ **do**
4:     Solve $P(H^q)$, and obtain its optimal solution $(\mathbf{x}^q, \mathbf{y}^q, \rho^q, \mathbf{w}^q, \mathbf{z}^q, \mathbf{U}^q)$.
5:     Update the lower bound: $LB^q \leftarrow v(P(H^q))$ using (22).
6:     Update the upper bound: $UB^q \leftarrow \min\{UB^{q-1}, Z(\mathbf{x}^q, \mathbf{y}^q)\}$ using (23).
7:     Generate a new set of points $U_j^{h_{new}} = \dfrac{\sum_{i \in I} \lambda_i x_{ij}^q}{\sum_{k \in K} \mu_{jk} y_{jk}^q - \sum_{i \in I} \lambda_i x_{ij}^q}$   $\forall j$ using (10).
8:     $H^{q+1} \leftarrow H^q \cup \{h_{new}\}$.
9:     $q \leftarrow q + 1$
10: **end while**

---

**Proposition 4:** *The constraint generation algorithm to solve [P(H)] is finite.*

**Proof:**

Since $x_{ij}, y_{jk} \in \{0, 1\}$ $\forall i, j, k$ and $U_j = \dfrac{\sum_{i \in I} \lambda_i x_{ij}}{\sum_{k \in K} \mu_{jk} y_{jk} - \sum_{i \in I} \lambda_i x_{ij}}$, $U_j$ can take only a finite set of values. Therefore, in order to prove the finiteness of Algorithm 1, it is sufficient to prove that the generated values of $U_j^h$ are not repeated.

Consider an iteration $q$, wherein the Algorithm 1 has not yet converged, that is $UB^q > LB^q$. Further, suppose $(\mathbf{x}^2, \mathbf{y}^2, \rho^2, \mathbf{w}^2, \mathbf{z}^2, \mathbf{U}^2)$ is a solution to $[P(H^q)]$. The new points $U_j^{h_{new}}$ generated at iteration $q$ are given by:

$$U_j^{h_{new}} = \frac{\sum_{i \in I} \lambda_i x_{ij}^q}{\sum_{k \in K} \mu_{jk} y_{jk}^q - \sum_{i \in I} \lambda_i x_{ij}^q} \quad \forall i, j, k$$

Suppose the values of $U_j^{h_{new}}$ are already generated at iteration $q^0 < q$ $\forall j \in J$. Then,

We have $U_j^{h_{new}}$ generated in iteration $q$ is given by : $U_j^{h_{new}} = \dfrac{\sum_{i \in I} \lambda_i x_{ij}^q}{\sum_{k \in K} \mu_{jk} y_{jk}^q - \sum_{i \in I} \lambda_i x_{ij}^q}$

$$(13) \Rightarrow \frac{U_j^{h_{new}}}{1 + U_j^{h_{new}}} \leq \frac{1}{1 + U_j^{h_{new}}} U_j^q + \left( \frac{U_j^{h_{new}}}{1 + U_j^{h_{new}}} \right)^2 \quad \forall j$$

$$\Rightarrow U_j^{h_{new}} \leq U_j^q \quad \forall j$$

We now have:

$$
\begin{aligned}
LB^q &= \sum_{j\in J}\sum_{k\in K} f_{jk}y^q_{jk} + \sum_{i\in I}\sum_{j\in J} c_{ij}x^q_{ij} + \frac{d}{2}\sum_{j\in J}\left\{ U^q_j + \rho^q_j + \sum_{k\in K} Cv^2_{jk}(w^q_{jk} - z^q_{jk})\right\} \\
&\geq \sum_{j\in J}\sum_{k\in K} f_{jk}y^q_{jk} + \sum_{i\in I}\sum_{j\in J} c_{ij}x^q_{ij} + \frac{d}{2}\sum_{j\in J}\left\{ U^{h_{new}}_j + \rho^q_j + \sum_{k\in K} Cv^2_{jk}(w^q_{jk} - z^q_{jk})\right\} \\
&= \sum_{j\in J}\sum_{k\in K} f_{jk}y^q_{jk} + \sum_{i\in I}\sum_{j\in J} c_{ij}x^q_{ij} + \frac{d}{2}\sum_{j\in J}\left\{ \frac{\sum_{i\in I}\lambda_i x^q_{ij}}{\sum_{k\in K}\mu_{jk}y^q_{jk} - \sum_{i\in I}\lambda_i x^q_{ij}} + \rho^q_j + \sum_{k\in K} Cv^2_{jk}(w^q_{jk} - z^q_{jk})\right\} \\
&= \sum_{k\in K} f_{jk}y^q_{jk} + \sum_{i\in I}\sum_{j\in J} c_{ij}x^q_{ij} + d\left\{ \frac{\left(1 + \sum_{k\in K} Cv^2_{jk}y_{jk}\right)\left(\sum_{i\in I}\lambda_i x_{ij}\right)^2}{2\sum_{k\in K}\mu_{jk}y_{jk}\left(\sum_{k\in K}\mu_{jk}y_{jk} - \sum_{i\in I}\lambda_i x_{ij}\right)} + \frac{\sum_{i\in I}\lambda_i x_{ij}}{\sum_{k\in K}\mu_{jk}y_{jk}}\right\} \\
&= Z(\mathbf{x}^q, \mathbf{y}^q) \\
&\geq \min\{UB^{q-1}, Z(\mathbf{x}^q, \mathbf{y}^q)\} \\
&= UB^q
\end{aligned}
$$

This contradicts are initial assumption $UB^q > LB^q$. Therefore, at any given iteration, $U^{h_{new}}_j$ will be different from the previously generated points at least for any one $j$. Furthermore, the number of values that $U^h_j$ can take is finite. Hence, the algorithm should terminate in a finite number of iterations.
∎

## 5. Computational Results

We present our computational experiments with the solution approach proposed in Section 4. The solution procedures are coded in C++ (Visual Studio 2010), while $[P(H^q)]$ at every iteration $q$ is solved using IBM ILOG CPLEX 12.4 on a personal laptop with Intel Core i5-3230M, 2.60 GHz CPU; 4 GB RAM; and Windows 64-bit operating system. The test instances are generated using a combination of the schemes used by Amiri (1998) and Holmberg et al. (1999), described in detail in Section 5.1. We generate four sets of test problems by varying the combination of the number of user nodes and the number of potential facility locations $(|I|, |J|)$ as (100, 10), (200, 15), (300, 20), and (400, 25). The number of capacity levels $(|K|)$ is set at 5, and the tolerance level for optimality $\epsilon$ is set at $10^{-5}$ in all the experiments.

### 5.1. Data Generation

The co-ordinates of the nodes representing user nodes are randomly generated using a uniform distribution $U \sim (10, 300)$. The mean demand rate at any user node $i$ is randomly generated as $\lambda_i \sim U(10, 50)$, following Holmberg et al. (1999). The locations of the potential facilities are obtained from the solution to a $p$-median facility location problem (Holmberg et al., 1999). The other parameters are generated as follows:

- *Capacity Levels* $(\mu_{jk})$ are set as: $\mu_{j1} = 0.50 * \mu_{j3}$; $\mu_{j2} = 0.75 * \mu_{j3}$; $\mu_{j4} = 1.25 * \mu_{j3}$; $\mu_{j5} = 1.50 * \mu_{j3}$, where $\mu_{j3} = (1.25 \sum_{i \in I} \lambda_i)/(|J| * LR)$. The Load Ratio (LR), which is defined as the average facility utilization if all the potential facilities are assigned capacity level 3 ($k = 3$) to satisfy 125% of the total demand in the network, is set at 0.60 (Amiri, 1998).

- *Fixed Costs* $(f_{jk})$ are generated as: $f_{j3} = FCR * E_{jj_0}$, where $E_{jj_0}$ is the Euclidean distance between the facility site $j$ and a fixed point $j_0 = (155, 155)$, which is the center of the box within which the coordinates of user nodes are generated. FCR, called the Fixed Cost Ratio, is a constant set at 40. $f_{j1} = 0.60 * f_{j3}$, $f_{j2} = 0.85 * f_{j3}$, $f_{j4} = 1.15 * f_{j3}$, $f_{j5} = 1.35 * f_{j3}$ (Amiri, 1998). The chosen values of fixed cost for different capacity levels exhibit both an underlying economy as well as diseconomy of scale. For example, for a 25% increase in capacity (corresponding to $\mu_{j4}$ over $\mu_{j3}$), the capacity cost increases only 15%. However, for the a 50% increase in capacity (corresponding to $\mu_{j5}$ over $\mu_{j3}$), the capacity cost increases 35%.

- *Service Costs* $(c_{ij})$ are generated as: $c_{ij} = 5 * E_{ij}$, where $E_{ij}$ is the Euclidean distance between user node $i$ and the facility site $j$ (Holmberg et al., 1999).

- *Unit Queueing Delay Cost* $(d)$ is assumed to be one of the values from the set $\{1, 10, 25, 50, 100, 250, 500, 1000, 5000\}$.

*5.2. Analysis of Computational Results*

Table 1 compares the performance of the proposed solution approach with a priori cuts versus without a priori cuts for test instances with 100 user nodes and 10 facilities. The table reports the number of iterations required (#Iter) and the computation time in seconds (CPU) for different values of coefficient of variation of service times ($Cv$) and the unit queuing delay cost ($d$). The table also reports the percentage reduction in computation time as a result of adding a priori cuts, which is computed as: (CPU time without a priori cuts - CPU time with a priori cuts)$\times 100$/(CPU time without a priori cuts). As described in Section 4.1, a priori cuts for the function $\rho(U) = U/(1 + U)$ are generated based on the piecewise linear approximation $\widehat{\rho}(U)$ of the function $\rho(U)$ such that the approximation error (measured by $\widehat{\rho}(U) - \rho(U)$) is at most 0.001 (Elhedhli, 2005). Figure 1 shows for $(|I| = 100, |J| = 10)$ the percentage reduction in the number of iterations and the CPU times as a result of addition of a priori cuts for different values of $d$ and $Cv$.

Following observations can be made from Table 1 and Figure 1. The results in Table 1 show that without a priori cuts, the problem takes, on an average, 949 seconds and 9 iterations to solve. The addition of a priori cuts reduces these numbers to less than 4 seconds of CPU time and less than 3 iterations. The percentage reduction in CPU time due to the addition of a priori cuts varies between 24.21% to 99.97%, with an average reduction

12

of 75.25%. Plots in Figure 1 further show that the benefits of a priori cuts, in terms of % reduction in CPU time and number of iterations, increase significantly with an increase in the unit queuing delay cost ($d$). This is because, as Table 1 suggests, without a priori cuts, the computation time and the number of iterations required by the proposed solution approach increases significantly with an increase in $d$. However, with the addition of a priori cuts, the computation time and the number of iterations required are almost *insensitive* to $d$.

Tables 2-5 summarize the computational results for four sets of instances ($|I| = 100$, $|J| = 10$; $|I| = 200$, $|J| = 15$; $|I| = 300$, $|J| = 20$; $|I| = 400$, $|J| = 25$), for different values of the unit queueing delay cost ($d = 1, 10, 25, 50, 100, 250, 500, 1000, 5000$) and coefficient of variation of service times ($Cv = 0, 1, 0.5, 1.5, 2, 2.5$). In solving each of the problem instances, we use a priori set of 32 cuts, which corresponds to a maximum approximation error of 0.001 in the linear approximation of $\rho(U)$. The tables report for each problem instance the total cost (TC); fixed cost (FC), access cost (AC), and delay cost (DC), expressed as a percentage of the total cost; computation time in seconds (CPU); number of iterations for convergence (#Iter); number of facilities opened (#Facility); and the minimum, maximum, and average facility utilizations among the open facilities. Figure 2 shows the effect of varying $d$ on FC, AC, DC, and TC for different values of $Cv$. Figure 3 shows the effect of varying $d$ on the maximum and average facility utilizations. Following observations can be made from the figures:

- An increase in $d$ or $Cv$ increases TC, as expected. An increase in $d$ also results in a higher DC, which is expected, provided the expected total number of users waiting at different service facilities in the network ($\sum_j \Lambda_j E[w_j]$) remains unchanged with an increase in $d$. However, an increase in $d$ implies a larger penalty for congestion (waiting customers), which forces the system to either attain more uniform utilization ($\rho_j = \Lambda_j/\mu_j$) among the open facilities or to install a larger total service capacity in the network (either by opening more facilities or by installing larger capacities at fewer, more or the same number of opened facilities). In either case, the maximum facility utilization decreases and the average facility utilization either decreases or remains constant, as evident from Figure 3. This results in a decrease in the expected total number of users waiting in the network ($\sum_j \Lambda_j E[w_j]$). We observe from our results that the percentage decrease in the expected total number of waiting users in the network is smaller compared to the percentage increase in $d$, such that the net effect is always an increase the total delay cost (DC).

- For a fixed network configuration (location and allocation of service facilities), an increase in $Cv$ is expected to increase the expected total number of waiting users in the network, and hence increase DC. However, when the location and allocation of service facilities are also decision variables, as they are in the current problem, Figure 2 suggests that an increase in $Cv$ may sometimes result in a decrease in DC,

13

which initially appears to be counter intuitive. However, this can be justified as follows. To counter the increase in congestion at a higher value of $Cv$, the system chooses to either attain more uniform utilization among the open facilities or to install a larger total service capacity in the system, thereby resulting in a decrease in the expected total number of waiting users in the network, and hence a decrease in DC.

- The fixed cost (FC) changes non-monotonically with an increase in $d$ or $Cv$, which also seems counter intuitive since to counter the effect of increase in $d$ or $Cv$, the system is expected to either attain more uniform utilization among the open facilities or to install a larger total service capacity in the system, neither of which should decrease FC. However, this seemingly counter intuitive observation can be justified as follows. In an attempt to equalize utilizations among open facilities, the system may choose to redistribute the total service capacity more uniformly among the open facilities, in which case some facilities may exhibit economies of scale while others may exhibit diseconomies of scale in fixed costs (as described in Section 5.1). This may result in either an increase or decrease in FC depending on the net effect of economies and diseconomies of scale. In addition, since the fixed cost of opening a service facility with a given capacity level ($f_{jk}$) depends on its distance from a fixed point ($j_0$), as described in Section 5.1, an increase in $d$ or $Cv$ may result in a decrease in FC if the system chooses to open most or all service facilities closer to $j_0$ (even though at the same or higher capacity levels) at the higher value of $d$ or $Cv$. In such a case, AC is likely to increase with an increase in $d$ or $Cv$. On the other hand, if the service facilities get more widely dispersed closer to user nodes in response to an increase in $d$ or $Cv$, then AC is likely to decrease.

- Comparison of CPU times across Tables 2-5 shows, as expected, that the computation time increases as the number of user nodes and potential facilities ($|I|, |J|$) increase.

- The proposed solution approach succeeds in finding exact (within an optimality gap of $10^{-5}$) solutions to the four sets of problem instances ($|I| = 100, |J| = 10; |I| = 200, |J| = 15; |I| = 300, |J| = 20; |I| = 400, |J| = 25$) within an average CPU time of 4, 29, 100 and 323 seconds, respectively, while the maximum CPU times for these sets are 17, 64, 668, and 1585 seconds, respectively. The average number of iterations are 2, 3, 3, and 3 whereas the maximum number of iterations are 3, 3, 3, and 6, respectively for the four sets of problem instances. This demonstrates the efficiency of the solution approach where the number of iterations/cuts imply that only a fraction of the exponential number constraints is required.

14

Figure 1: Effect of Adding a-priori Cuts on the Number of Iterations and the CPU Times of the Algorithm



Figure 2: Effect of Varying Unit Delay Cost on the Fixed Cost, Access Cost, Delay Cost, and the Total Cost

Figure 3: Effect of Varying Unit Delay Cost on Maximum and Average Facility Utilization

Table 1: Effect of Adding a-priori Cuts on the Performance of the Solution Method: Instances with 100 User Nodes and 10 Potential Facilities

| | | Without a-priori cuts | | With a-priori cuts | | % Reduction |
|---|---|---|---|---|---|---|
| $Cv$ | $d$ | # Iter | CPU | # Iter | CPU | in CPU times |
| 0 | 1 | 3 | 2.2 | 2 | 1.6 | 29 |
| | 10 | 6 | 8.0 | 2 | 1.8 | 77 |
| | 25 | 6 | 5.4 | 3 | 3.0 | 45 |
| | 50 | 7 | 12.1 | 2 | 3.7 | 69 |
| | 100 | 7 | 13.5 | 2 | 3.4 | 75 |
| | 250 | 8 | 11.0 | 2 | 2.0 | 81 |
| | 500 | 8 | 12.1 | 2 | 1.6 | 87 |
| | 1000 | 12 | 21.7 | 3 | 2.5 | 88 |
| | 5000 | 20 | 24.8 | 2 | 2.5 | 90 |
| | | | | | | |
| 0.5 | 1 | 3 | 2.1 | 2 | 1.6 | 24 |
| | 10 | 6 | 5.9 | 3 | 3.6 | 38 |
| | 25 | 6 | 6.1 | 2 | 2.4 | 61 |
| | 50 | 7 | 10.5 | 2 | 3.5 | 67 |
| | 100 | 7 | 12.3 | 2 | 4.0 | 67 |
| | 250 | 11 | 13.8 | 2 | 2.3 | 83 |
| | 500 | 12 | 19.1 | 3 | 2.9 | 85 |
| | 1000 | 14 | 22.1 | 3 | 2.3 | 90 |
| | 5000 | 18 | 28.0 | 2 | 1.8 | 94 |
| | | | | | | |
| 1 | 1 | 4 | 3.2 | 2 | 1.3 | 59 |
| | 10 | 6 | 7.1 | 2 | 2.2 | 68 |
| | 25 | 6 | 8.7 | 2 | 3.8 | 56 |
| | 50 | 8 | 16.0 | 2 | 4.1 | 74 |
| | 100 | 7 | 12.1 | 2 | 3.3 | 72 |
| | 250 | 9 | 16.7 | 2 | 1.3 | 92 |
| | 500 | 16 | 22.2 | 3 | 2.4 | 89 |
| | 1000 | 13 | 24.4 | 2 | 2.3 | 91 |
| | 5000 | 15 | 38.7 | 3 | 8.3 | 79 |
| | | | | | | |
| 1.5 | 1 | 4 | 2.6 | 2 | 1.6 | 39 |
| | 10 | 6 | 6.4 | 2 | 2.3 | 64 |
| | 25 | 8 | 14.7 | 2 | 3.2 | 78 |
| | 50 | 7 | 14.8 | 2 | 4.2 | 72 |
| | 100 | 8 | 17.1 | 2 | 2.0 | 88 |
| | 250 | 13 | 34.2 | 2 | 1.8 | 95 |
| | 500 | 13 | 25.7 | 2 | 2.2 | 92 |
| | 1000 | 11 | 29.3 | 2 | 3.2 | 89 |
| | 5000 | 12 | 60.6 | 3 | 16.8 | 72 |
| | | | | | | |
| 2 | 1 | 4 | 2.8 | 2 | 1.4 | 49 |
| | 10 | 6 | 9.1 | 2 | 3.4 | 63 |
| | 25 | 6 | 11.7 | 2 | 4.4 | 62 |
| | 50 | 8 | 15.7 | 3 | 3.9 | 75 |
| | 100 | 9 | 19.8 | 2 | 1.4 | 93 |
| | 250 | 12 | 25.3 | 3 | 2.3 | 91 |
| | 500 | 12 | 22.4 | 2 | 2.9 | 87 |
| | 1000 | 12 | 38.5 | 3 | 4.3 | 89 |
| | 5000 | 12 | 162.5 | 3 | 15.6 | 90 |
| | | | | | | |
| 2.5 | 1 | 6 | 5.2 | 2 | 1.7 | 67 |
| | 10 | 7 | 12.2 | 3 | 4.7 | 62 |
| | 25 | 7 | 12.2 | 2 | 3.9 | 68 |
| | 50 | 10 | 23.5 | 2 | 2.5 | 89 |
| | 100 | 12 | 31.3 | 2 | 2.0 | 93 |
| | 250 | 10 | 25.4 | 2 | 2.3 | 91 |
| | 500 | 11 | 26.8 | 3 | 4.9 | 82 |
| | 1000 | 11 | 50.0 | 3 | 4.8 | 90 |
| | 5000 | 14 | 50154.7 | 3 | 13.1 | 100 |
| | | | | | | |
| | Min. | 3 | 2.1 | 2 | 1.3 | 24 |
| | Avg. | 9 | 948.8 | 2 | 3.6 | 75 |
| | Max. | 20 | 50154.7 | 3 | 16.8 | 100 |

17

Table 2: Computational Performance: Instances with 100 User Nodes and 10 Potential Facilities

| $Cv$ | $d$ | $TC$ | FC (%) | AC (%) | DC (%) | # Facility | % Utilization | | | # Iter | CPU |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | Min. | Max. | Avg. | | |
| 0 | 1 | 42,405 | 47 | 53 | 0 | 7 | 89 | 98 | 67 | 2 | 2 |
| | 10 | 43,248 | 46 | 52 | 2 | 7 | 95 | 97 | 67 | 2 | 2 |
| | 25 | 43,962 | 46 | 51 | 2 | 7 | 89 | 94 | 64 | 3 | 3 |
| | 50 | 44,955 | 47 | 50 | 3 | 7 | 83 | 93 | 61 | 2 | 4 |
| | 100 | 46,308 | 47 | 48 | 5 | 7 | 81 | 89 | 59 | 2 | 3 |
| | 250 | 48,741 | 40 | 53 | 7 | 6 | 74 | 85 | 47 | 2 | 2 |
| | 500 | 51,438 | 34 | 56 | 10 | 5 | 72 | 81 | 38 | 2 | 2 |
| | 1,000 | 55,569 | 36 | 51 | 13 | 5 | 65 | 72 | 34 | 3 | 3 |
| | 5,000 | 79,578 | 36 | 32 | 32 | 6 | 48 | 56 | 32 | 2 | 2 |
| | | | | | | | | | | | |
| 0.5 | 1 | 42,429 | 47 | 53 | 0 | 7 | 89 | 98 | 67 | 2 | 2 |
| | 10 | 43,392 | 47 | 52 | 1 | 7 | 89 | 97 | 64 | 3 | 4 |
| | 25 | 44,186 | 46 | 51 | 3 | 7 | 89 | 94 | 64 | 2 | 2 |
| | 50 | 45,251 | 46 | 50 | 4 | 7 | 83 | 93 | 61 | 2 | 3 |
| | 100 | 46,720 | 46 | 48 | 6 | 7 | 81 | 89 | 59 | 2 | 4 |
| | 250 | 49,255 | 35 | 58 | 7 | 5 | 75 | 82 | 40 | 2 | 2 |
| | 500 | 52,233 | 37 | 54 | 9 | 5 | 69 | 73 | 35 | 3 | 3 |
| | 1,000 | 56,507 | 35 | 50 | 14 | 5 | 65 | 72 | 34 | 3 | 2 |
| | 5,000 | 81,897 | 35 | 31 | 34 | 6 | 48 | 56 | 32 | 2 | 2 |
| | | | | | | | | | | | |
| 1 | 1 | 42,501 | 47 | 53 | 0 | 7 | 89 | 98 | 67 | 2 | 1 |
| | 10 | 43,684 | 47 | 51 | 2 | 7 | 89 | 94 | 64 | 2 | 2 |
| | 25 | 44,802 | 47 | 50 | 3 | 7 | 83 | 93 | 61 | 2 | 4 |
| | 50 | 46,013 | 47 | 49 | 4 | 7 | 81 | 89 | 59 | 2 | 4 |
| | 100 | 47,555 | 41 | 54 | 5 | 6 | 74 | 85 | 47 | 2 | 3 |
| | 250 | 50,479 | 35 | 57 | 8 | 5 | 72 | 81 | 38 | 2 | 1 |
| | 500 | 53,859 | 37 | 53 | 10 | 5 | 65 | 70 | 34 | 3 | 2 |
| | 1,000 | 59,040 | 35 | 48 | 17 | 5 | 61 | 70 | 33 | 2 | 2 |
| | 5,000 | 88,546 | 41 | 25 | 34 | 7 | 42 | 52 | 32 | 3 | 8 |
| | | | | | | | | | | | |
| 1.5 | 1 | 42,615 | 47 | 52 | 1 | 7 | 95 | 98 | 67 | 2 | 2 |
| | 10 | 44,135 | 46 | 51 | 3 | 7 | 89 | 94 | 64 | 2 | 2 |
| | 25 | 45,525 | 46 | 50 | 4 | 7 | 83 | 92 | 61 | 2 | 3 |
| | 50 | 46,920 | 47 | 48 | 5 | 7 | 74 | 85 | 56 | 2 | 4 |
| | 100 | 48,671 | 41 | 53 | 6 | 6 | 72 | 81 | 46 | 2 | 2 |
| | 250 | 52,162 | 37 | 54 | 8 | 5 | 69 | 73 | 35 | 2 | 2 |
| | 500 | 56,163 | 37 | 50 | 12 | 5 | 61 | 70 | 33 | 2 | 2 |
| | 1,000 | 62,778 | 35 | 45 | 20 | 5 | 62 | 66 | 32 | 2 | 3 |
| | 5,000 | 97,044 | 38 | 23 | 39 | 7 | 42 | 48 | 32 | 3 | 17 |
| | | | | | | | | | | | |
| 2 | 1 | 42,763 | 47 | 52 | 1 | 7 | 95 | 97 | 67 | 2 | 1 |
| | 10 | 44,711 | 47 | 50 | 3 | 7 | 83 | 93 | 61 | 2 | 3 |
| | 25 | 46,278 | 47 | 48 | 5 | 7 | 81 | 89 | 59 | 2 | 4 |
| | 50 | 47,791 | 41 | 54 | 5 | 6 | 74 | 85 | 47 | 3 | 4 |
| | 100 | 49,903 | 35 | 57 | 7 | 5 | 72 | 81 | 38 | 2 | 1 |
| | 250 | 53,941 | 37 | 53 | 10 | 5 | 65 | 70 | 34 | 3 | 2 |
| | 500 | 59,023 | 35 | 48 | 17 | 5 | 61 | 70 | 33 | 2 | 3 |
| | 1,000 | 66,506 | 42 | 38 | 20 | 6 | 49 | 59 | 33 | 3 | 4 |
| | 5,000 | 107,202 | 41 | 19 | 40 | 8 | 36 | 44 | 32 | 3 | 16 |
| | | | | | | | | | | | |
| 2.5 | 1 | 42,954 | 47 | 52 | 1 | 7 | 95 | 97 | 67 | 2 | 2 |
| | 10 | 45,238 | 46 | 50 | 4 | 7 | 83 | 93 | 61 | 3 | 5 |
| | 25 | 46,993 | 47 | 48 | 5 | 7 | 74 | 85 | 56 | 2 | 4 |
| | 50 | 48,735 | 41 | 53 | 6 | 6 | 72 | 81 | 46 | 2 | 3 |
| | 100 | 51,249 | 38 | 55 | 7 | 5 | 69 | 73 | 35 | 2 | 2 |
| | 250 | 55,950 | 37 | 51 | 12 | 5 | 61 | 70 | 33 | 2 | 2 |
| | 500 | 61,983 | 43 | 41 | 16 | 6 | 48 | 60 | 34 | 3 | 5 |
| | 1,000 | 70,952 | 40 | 36 | 23 | 6 | 48 | 56 | 32 | 3 | 5 |
| | 5,000 | 118,106 | 44 | 15 | 41 | 9 | 34 | 39 | 32 | 3 | 13 |
| | Min. | 42405 | 34 | 15 | 0 | 5 | 34 | 39 | 32 | 2 | 1 |
| | Avg. | 54645 | 42 | 48 | 10 | 6 | 72 | 80 | 48 | 2 | 4 |
| | Max. | 118106 | 47 | 58 | 41 | 9 | 95 | 98 | 67 | 3 | 17 |

Table 3: Computational Performance: Instances with 200 User Nodes and 15 Potential Facilities

| Cv | d | TC | FC (%) | AC (%) | DC (%) | # Facility | % Utilization | | | # Iter | CPU |
|----|----|------|--------|--------|--------|-----------|------|------|------|--------|-----|
| | | | | | | | Min. | Max. | Avg. | | |
| 0 | 1 | 64,897 | 55 | 44 | 0 | 14 | 91 | 99 | 90 | 2 | 18 |
| | 10 | 66,409 | 55 | 43 | 2 | 14 | 76 | 98 | 87 | 2 | 18 |
| | 25 | 67,976 | 54 | 43 | 3 | 14 | 78 | 95 | 84 | 2 | 23 |
| | 50 | 69,594 | 55 | 41 | 4 | 14 | 71 | 93 | 80 | 2 | 20 |
| | 100 | 72,116 | 49 | 46 | 5 | 12 | 67 | 90 | 66 | 3 | 60 |
| | 250 | 76,189 | 42 | 51 | 7 | 9 | 62 | 82 | 47 | 2 | 26 |
| | 500 | 80,387 | 42 | 49 | 9 | 9 | 56 | 76 | 42 | 2 | 17 |
| | 1,000 | 86,593 | 43 | 45 | 12 | 9 | 45 | 72 | 37 | 2 | 14 |
| | 5,000 | 120,794 | 36 | 33 | 32 | 9 | 48 | 60 | 32 | 3 | 28 |
| | | | | | | | | | | | |
| 0.5 | 1 | 64,961 | 55 | 44 | 1 | 14 | 91 | 99 | 90 | 2 | 16 |
| | 10 | 66,685 | 55 | 43 | 2 | 14 | 88 | 96 | 87 | 2 | 22 |
| | 25 | 68,372 | 55 | 42 | 3 | 14 | 76 | 94 | 82 | 2 | 23 |
| | 50 | 70,108 | 55 | 41 | 5 | 14 | 71 | 93 | 80 | 2 | 19 |
| | 100 | 72,765 | 49 | 45 | 6 | 12 | 70 | 87 | 66 | 2 | 31 |
| | 250 | 76,982 | 43 | 51 | 7 | 9 | 62 | 82 | 45 | 3 | 40 |
| | 500 | 81,364 | 43 | 48 | 9 | 9 | 50 | 76 | 40 | 2 | 15 |
| | 1,000 | 87,860 | 42 | 44 | 13 | 9 | 50 | 72 | 37 | 3 | 10 |
| | 5,000 | 124,123 | 40 | 29 | 31 | 10 | 41 | 55 | 32 | 5 | 57 |
| | | | | | | | | | | | |
| 1 | 1 | 65,136 | 55 | 44 | 1 | 14 | 91 | 99 | 90 | 2 | 19 |
| | 10 | 67,401 | 55 | 42 | 3 | 14 | 88 | 96 | 87 | 4 | 54 |
| | 25 | 69,295 | 55 | 41 | 3 | 14 | 71 | 93 | 80 | 2 | 19 |
| | 50 | 71,555 | 54 | 41 | 5 | 14 | 75 | 89 | 77 | 3 | 48 |
| | 100 | 74,309 | 45 | 50 | 5 | 10 | 67 | 83 | 52 | 2 | 32 |
| | 250 | 78,810 | 43 | 50 | 7 | 9 | 56 | 76 | 42 | 2 | 23 |
| | 500 | 83,774 | 44 | 46 | 10 | 9 | 46 | 72 | 38 | 2 | 17 |
| | 1,000 | 91,294 | 41 | 43 | 15 | 9 | 46 | 68 | 36 | 3 | 17 |
| | 5,000 | 132,533 | 38 | 27 | 35 | 10 | 43 | 51 | 32 | 4 | 35 |
| | | | | | | | | | | | |
| 1.5 | 1 | 65,396 | 55 | 44 | 1 | 14 | 92 | 98 | 90 | 2 | 22 |
| | 10 | 68,258 | 55 | 42 | 3 | 14 | 76 | 94 | 82 | 2 | 23 |
| | 25 | 70,557 | 54 | 41 | 5 | 14 | 75 | 93 | 79 | 2 | 24 |
| | 50 | 73,197 | 49 | 46 | 5 | 12 | 67 | 87 | 64 | 3 | 59 |
| | 100 | 76,017 | 45 | 48 | 6 | 10 | 64 | 82 | 51 | 3 | 35 |
| | 250 | 81,115 | 43 | 48 | 9 | 9 | 51 | 76 | 40 | 2 | 20 |
| | 500 | 86,919 | 43 | 45 | 12 | 9 | 50 | 70 | 37 | 3 | 17 |
| | 1,000 | 95,844 | 42 | 41 | 17 | 9 | 45 | 65 | 34 | 3 | 23 |
| | 5,000 | 144,737 | 39 | 24 | 38 | 11 | 38 | 48 | 32 | 4 | 53 |
| | | | | | | | | | | | |
| 2 | 1 | 65,675 | 56 | 43 | 1 | 14 | 76 | 98 | 87 | 2 | 16 |
| | 10 | 69,116 | 56 | 41 | 3 | 14 | 71 | 93 | 80 | 2 | 21 |
| | 25 | 72,011 | 54 | 41 | 6 | 14 | 75 | 89 | 77 | 2 | 36 |
| | 50 | 74,649 | 45 | 50 | 5 | 10 | 67 | 83 | 52 | 3 | 42 |
| | 100 | 77,852 | 43 | 50 | 7 | 9 | 56 | 79 | 43 | 2 | 23 |
| | 250 | 83,654 | 44 | 46 | 9 | 9 | 50 | 72 | 37 | 3 | 15 |
| | 500 | 90,632 | 43 | 44 | 13 | 9 | 45 | 67 | 35 | 3 | 22 |
| | 1,000 | 101,278 | 46 | 36 | 18 | 10 | 41 | 56 | 34 | 5 | 52 |
| | 5,000 | 159,170 | 39 | 20 | 40 | 12 | 37 | 44 | 32 | 4 | 43 |
| | | | | | | | | | | | |
| 2.5 | 1 | 65,957 | 56 | 43 | 1 | 14 | 76 | 98 | 87 | 2 | 20 |
| | 10 | 70,031 | 55 | 41 | 4 | 14 | 73 | 93 | 80 | 2 | 20 |
| | 25 | 73,347 | 49 | 45 | 5 | 12 | 70 | 86 | 64 | 3 | 53 |
| | 50 | 76,108 | 45 | 48 | 6 | 10 | 64 | 80 | 51 | 3 | 40 |
| | 100 | 79,628 | 44 | 49 | 7 | 9 | 50 | 76 | 40 | 2 | 24 |
| | 250 | 86,404 | 44 | 45 | 11 | 9 | 50 | 68 | 37 | 2 | 16 |
| | 500 | 94,597 | 42 | 42 | 16 | 9 | 49 | 65 | 34 | 4 | 32 |
| | 1,000 | 106,997 | 45 | 34 | 22 | 10 | 44 | 55 | 33 | 3 | 34 |
| | 5,000 | 175,392 | 40 | 17 | 43 | 13 | 33 | 41 | 32 | 6 | 64 |
| | Min. | 64,897 | 36 | 17 | 0 | 9 | 33 | 41 | 32 | 2 | 10 |
| | Avg. | 84,015 | 48 | 42 | 10 | 11 | 62 | 80 | 57 | 3 | 29 |
| | Max. | 175,392 | 56 | 51 | 43 | 14 | 92 | 99 | 90 | 6 | 64 |

Table 4: Computational Performance: Instances with 300 User Nodes and 20 Potential Facilities

| $Cv$ | $d$ | $TC$ | FC (%) | AC (%) | DC (%) | # Facility | % Utilization | | | # Iter | CPU |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | Min. | Max. | Avg. | | |
| 0 | 1 | 87,130 | 57 | 43 | 1 | 18 | 84 | 99 | 87 | 2 | 668 |
| | 10 | 89,241 | 56 | 42 | 2 | 18 | 91 | 97 | 84 | 3 | 52 |
| | 25 | 91,380 | 55 | 41 | 3 | 18 | 83 | 96 | 83 | 3 | 56 |
| | 50 | 93,915 | 54 | 41 | 5 | 18 | 85 | 92 | 81 | 3 | 71 |
| | 100 | 97,458 | 55 | 39 | 6 | 18 | 76 | 90 | 77 | 3 | 220 |
| | 250 | 103,111 | 47 | 45 | 8 | 14 | 61 | 84 | 54 | 2 | 75 |
| | 500 | 109,103 | 45 | 45 | 10 | 13 | 47 | 76 | 46 | 4 | 119 |
| | 1,000 | 117,835 | 42 | 46 | 12 | 11 | 41 | 72 | 36 | 3 | 72 |
| | 5,000 | 163,701 | 36 | 33 | 31 | 12 | 42 | 61 | 32 | 3 | 93 |
| | | | | | | | | | | | |
| 0.5 | 1 | 87,242 | 57 | 43 | 1 | 18 | 86 | 99 | 87 | 2 | 262 |
| | 10 | 89,583 | 56 | 42 | 2 | 18 | 91 | 97 | 84 | 2 | 38 |
| | 25 | 91,976 | 55 | 41 | 4 | 18 | 83 | 95 | 82 | 2 | 35 |
| | 50 | 94,704 | 55 | 41 | 5 | 18 | 76 | 92 | 79 | 2 | 48 |
| | 100 | 98,312 | 49 | 45 | 5 | 15 | 76 | 88 | 62 | 2 | 91 |
| | 250 | 104,193 | 46 | 47 | 7 | 13 | 56 | 79 | 48 | 3 | 91 |
| | 500 | 110,477 | 45 | 45 | 10 | 13 | 47 | 75 | 44 | 2 | 50 |
| | 1,000 | 119,613 | 41 | 46 | 13 | 11 | 46 | 71 | 36 | 4 | 93 |
| | 5,000 | 167,997 | 39 | 29 | 32 | 13 | 40 | 57 | 32 | 3 | 66 |
| | | | | | | | | | | | |
| 1 | 1 | 87,523 | 56 | 43 | 1 | 18 | 91 | 99 | 87 | 3 | 114 |
| | 10 | 90,593 | 55 | 42 | 3 | 18 | 91 | 96 | 84 | 2 | 59 |
| | 25 | 93,513 | 55 | 41 | 4 | 18 | 85 | 92 | 81 | 3 | 70 |
| | 50 | 96,690 | 55 | 39 | 6 | 18 | 76 | 90 | 77 | 3 | 119 |
| | 100 | 100,348 | 48 | 46 | 6 | 14 | 61 | 85 | 56 | 4 | 128 |
| | 250 | 106,782 | 45 | 46 | 8 | 13 | 51 | 78 | 47 | 4 | 98 |
| | 500 | 113,941 | 43 | 46 | 11 | 12 | 46 | 74 | 40 | 3 | 76 |
| | 1,000 | 124,301 | 40 | 45 | 14 | 11 | 41 | 67 | 34 | 2 | 47 |
| | 5,000 | 179,921 | 37 | 28 | 35 | 13 | 44 | 52 | 32 | 5 | 110 |
| | | | | | | | | | | | |
| 1.5 | 1 | 87,845 | 56 | 43 | 1 | 18 | 92 | 98 | 87 | 2 | 56 |
| | 10 | 91,847 | 55 | 41 | 4 | 18 | 83 | 95 | 82 | 2 | 38 |
| | 25 | 95,403 | 55 | 40 | 5 | 18 | 76 | 90 | 77 | 2 | 71 |
| | 50 | 98,800 | 50 | 45 | 5 | 15 | 61 | 84 | 60 | 3 | 115 |
| | 100 | 102,896 | 48 | 46 | 6 | 14 | 56 | 81 | 52 | 3 | 108 |
| | 250 | 110,165 | 45 | 45 | 9 | 13 | 51 | 74 | 44 | 3 | 71 |
| | 500 | 118,582 | 42 | 46 | 12 | 11 | 46 | 69 | 35 | 3 | 78 |
| | 1,000 | 130,498 | 42 | 41 | 17 | 12 | 42 | 64 | 34 | 4 | 120 |
| | 5,000 | 197,199 | 41 | 22 | 36 | 15 | 38 | 49 | 32 | 4 | 112 |
| | | | | | | | | | | | |
| 2 | 1 | 88,279 | 56 | 43 | 1 | 18 | 92 | 98 | 86 | 2 | 123 |
| | 10 | 93,271 | 55 | 41 | 4 | 18 | 85 | 92 | 81 | 3 | 74 |
| | 25 | 97,338 | 56 | 39 | 5 | 18 | 68 | 87 | 74 | 4 | 202 |
| | 50 | 100,919 | 49 | 44 | 6 | 15 | 68 | 83 | 59 | 3 | 107 |
| | 100 | 105,374 | 46 | 47 | 7 | 13 | 51 | 78 | 47 | 3 | 79 |
| | 250 | 113,993 | 44 | 46 | 10 | 12 | 43 | 70 | 38 | 2 | 45 |
| | 500 | 123,494 | 44 | 43 | 13 | 12 | 44 | 64 | 35 | 3 | 74 |
| | 1,000 | 137,761 | 46 | 36 | 18 | 13 | 40 | 59 | 33 | 3 | 77 |
| | 5,000 | 216,365 | 44 | 19 | 37 | 17 | 34 | 42 | 32 | 4 | 88 |
| | | | | | | | | | | | |
| 2.5 | 1 | 88,695 | 56 | 42 | 1 | 18 | 86 | 97 | 84 | 2 | 38 |
| | 10 | 94,674 | 55 | 41 | 5 | 18 | 76 | 92 | 79 | 3 | 83 |
| | 25 | 98,978 | 50 | 45 | 5 | 15 | 61 | 84 | 60 | 2 | 73 |
| | 50 | 102,958 | 48 | 46 | 6 | 14 | 56 | 80 | 51 | 2 | 71 |
| | 100 | 107,960 | 46 | 46 | 8 | 13 | 47 | 74 | 44 | 3 | 76 |
| | 250 | 117,909 | 46 | 43 | 12 | 13 | 44 | 68 | 39 | 2 | 51 |
| | 500 | 128,813 | 44 | 41 | 15 | 12 | 42 | 64 | 33 | 3 | 87 |
| | 1,000 | 145,395 | 45 | 34 | 21 | 13 | 42 | 57 | 32 | 5 | 116 |
| | 5,000 | 238,223 | 44 | 16 | 40 | 18 | 33 | 38 | 32 | 6 | 150 |
| | Min. | 87,130 | 36 | 16 | 1 | 11 | 33 | 38 | 32 | 2 | 35 |
| | Avg. | 113,782 | 49 | 41 | 10 | 15 | 62 | 79 | 58 | 3 | 100 |
| | Max. | 238,223 | 57 | 47 | 40 | 18 | 92 | 99 | 87 | 6 | 668 |

Table 5: Computational Performance: Instances with 400 User Nodes and 25 Potential Facilities

| Cv | d | TC | FC (%) | AC (%) | DC (%) | # Facility | % Utilization | | | # Iter | CPU |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | Min. | Max. | Avg. | | |
| 0 | 1 | 108,302 | 53 | 46 | 0 | 19 | 97 | 99 | 74 | 3 | 807 |
| | 10 | 110,755 | 53 | 45 | 2 | 19 | 90 | 97 | 72 | 2 | 70 |
| | 25 | 113,050 | 52 | 45 | 3 | 19 | 87 | 97 | 70 | 2 | 110 |
| | 50 | 115,745 | 52 | 44 | 4 | 19 | 78 | 94 | 67 | 3 | 241 |
| | 100 | 119,588 | 53 | 42 | 6 | 19 | 73 | 92 | 65 | 2 | 396 |
| | 250 | 126,660 | 50 | 43 | 7 | 17 | 66 | 85 | 52 | 2 | 332 |
| | 500 | 134,432 | 45 | 45 | 10 | 15 | 53 | 80 | 43 | 3 | 206 |
| | 1,000 | 145,174 | 43 | 44 | 13 | 14 | 48 | 76 | 37 | 3 | 106 |
| | 5,000 | 203,766 | 38 | 30 | 32 | 15 | 45 | 64 | 32 | 4 | 242 |
| | | | | | | | | | | | |
| 0.5 | 1 | 108,430 | 53 | 46 | 1 | 19 | 97 | 99 | 74 | 3 | 988 |
| | 10 | 111,135 | 53 | 45 | 2 | 19 | 92 | 97 | 71 | 3 | 98 |
| | 25 | 113,694 | 53 | 44 | 3 | 19 | 82 | 96 | 69 | 3 | 150 |
| | 50 | 116,659 | 53 | 43 | 4 | 19 | 78 | 92 | 66 | 4 | 1,033 |
| | 100 | 120,705 | 52 | 43 | 6 | 18 | 78 | 90 | 60 | 2 | 511 |
| | 250 | 128,061 | 50 | 42 | 7 | 17 | 62 | 83 | 51 | 3 | 237 |
| | 500 | 136,165 | 45 | 45 | 10 | 15 | 54 | 79 | 42 | 4 | 192 |
| | 1,000 | 147,499 | 42 | 44 | 14 | 14 | 48 | 76 | 36 | 2 | 65 |
| | 5,000 | 209,065 | 41 | 27 | 31 | 16 | 40 | 56 | 32 | 4 | 149 |
| | | | | | | | | | | | |
| 1 | 1 | 108,768 | 53 | 46 | 1 | 19 | 86 | 98 | 73 | 2 | 558 |
| | 10 | 112,176 | 52 | 45 | 2 | 19 | 81 | 97 | 70 | 2 | 75 |
| | 25 | 115,324 | 52 | 44 | 4 | 19 | 78 | 94 | 67 | 3 | 229 |
| | 50 | 118,776 | 53 | 42 | 5 | 19 | 73 | 92 | 65 | 2 | 396 |
| | 100 | 123,345 | 50 | 44 | 6 | 17 | 68 | 87 | 54 | 4 | 945 |
| | 250 | 131,521 | 50 | 41 | 9 | 17 | 58 | 81 | 49 | 3 | 183 |
| | 500 | 140,513 | 46 | 43 | 11 | 15 | 48 | 75 | 39 | 3 | 112 |
| | 1,000 | 153,517 | 43 | 42 | 15 | 14 | 46 | 71 | 34 | 3 | 131 |
| | 5,000 | 224,302 | 41 | 25 | 34 | 17 | 39 | 53 | 32 | 5 | 335 |
| | | | | | | | | | | | |
| 1.5 | 1 | 109,115 | 53 | 46 | 1 | 19 | 86 | 98 | 73 | 2 | 93 |
| | 10 | 113,548 | 53 | 44 | 3 | 19 | 84 | 94 | 69 | 3 | 151 |
| | 25 | 117,369 | 54 | 42 | 4 | 19 | 73 | 92 | 65 | 4 | 707 |
| | 50 | 121,414 | 51 | 43 | 6 | 18 | 73 | 88 | 59 | 3 | 734 |
| | 100 | 126,409 | 50 | 43 | 7 | 17 | 62 | 83 | 51 | 4 | 335 |
| | 250 | 135,945 | 46 | 45 | 9 | 15 | 54 | 78 | 41 | 3 | 181 |
| | 500 | 146,324 | 44 | 44 | 12 | 14 | 48 | 74 | 35 | 5 | 208 |
| | 1,000 | 161,516 | 42 | 40 | 17 | 14 | 45 | 68 | 33 | 3 | 113 |
| | 5,000 | 245,741 | 43 | 21 | 36 | 19 | 37 | 51 | 32 | 4 | 175 |
| | | | | | | | | | | | |
| 2 | 1 | 109,592 | 53 | 46 | 1 | 19 | 88 | 98 | 73 | 2 | 100 |
| | 10 | 115,071 | 52 | 45 | 3 | 19 | 78 | 94 | 68 | 3 | 218 |
| | 25 | 119,552 | 53 | 42 | 5 | 19 | 78 | 88 | 64 | 3 | 771 |
| | 50 | 124,030 | 51 | 44 | 6 | 17 | 69 | 85 | 53 | 2 | 350 |
| | 100 | 129,675 | 50 | 42 | 8 | 17 | 58 | 80 | 49 | 3 | 162 |
| | 250 | 140,567 | 46 | 43 | 10 | 15 | 48 | 74 | 39 | 3 | 114 |
| | 500 | 152,752 | 44 | 42 | 13 | 14 | 46 | 69 | 33 | 5 | 240 |
| | 1,000 | 170,954 | 48 | 33 | 19 | 16 | 43 | 58 | 33 | 4 | 177 |
| | 5,000 | 269,628 | 42 | 18 | 40 | 20 | 35 | 45 | 32 | 4 | 122 |
| | | | | | | | | | | | |
| 2.5 | 1 | 110,130 | 53 | 45 | 1 | 19 | 82 | 98 | 72 | 2 | 181 |
| | 10 | 116,603 | 54 | 43 | 4 | 19 | 78 | 93 | 66 | 2 | 373 |
| | 25 | 121,739 | 51 | 43 | 5 | 18 | 73 | 88 | 58 | 5 | 1,585 |
| | 50 | 126,527 | 51 | 43 | 6 | 17 | 62 | 83 | 51 | 2 | 97 |
| | 100 | 133,103 | 50 | 41 | 9 | 17 | 57 | 78 | 47 | 2 | 82 |
| | 250 | 145,560 | 46 | 42 | 12 | 15 | 47 | 71 | 37 | 3 | 133 |
| | 500 | 159,695 | 46 | 39 | 15 | 15 | 45 | 66 | 34 | 6 | 369 |
| | 1,000 | 181,019 | 46 | 32 | 22 | 16 | 43 | 56 | 33 | 3 | 146 |
| | 5,000 | 298,560 | 43 | 16 | 41 | 22 | 32 | 40 | 32 | 5 | 621 |
| | Min. | 108,302 | 38 | 16 | 0 | 14 | 32 | 40 | 32 | 2 | 65 |
| | Avg. | 140,727 | 49 | 41 | 10 | 17 | 64 | 81 | 52 | 3 | 323 |
| | Max. | 298,560 | 54 | 46 | 41 | 22 | 97 | 99 | 74 | 6 | 1,585 |

## 6. Conclusions

In this paper, we presented a class of location-allocation problems with immobile servers, stochastic demand and congestion. The model captures the trade-off among the fixed cost

21

of opening service facilities and equipping them with sufficient capacities, the access cost associated with users' travel to service facilities, and the queueing delay cost associated with customers waiting for service. Under the assumption that the customer demands follow a Poisson process and service times follow a general distribution, the facilities were modeled as a network of independent $M/G/1$ queues, whose locations, capacity levels and workload allocations are decision variables. We presented a non-linear IP formulation and a constraint generation based exact method to solve its linear MIP reformulation. The computational results indicate that the proposed approach provides optimal solution for problem instances of the size up to 400 nodes and 25 potential facility locations within reasonable computation times. Future research directions may include extending the proposed solution procedures to deal with systems with multiple servers and general demand processes.

## 7. Acknowledgements

## References

Aboolian, R., Berman, O., Drezner, Z., 2008. Location and allocation of service units on a congested network. IIE Transactions 40, 422–433.

Amiri, A., 1997. Solution procedures for the service system design problem. Computers and Operations Research 24, 49–60.

Amiri, A., 1998. The design of service systems with queuing time cost, workload capacities, and backup service. European Journal of Operational Research 104, 201–217.

Amiri, A., 2001. The multi-hour service system design problem. European Journal of Operational Research 128, 625–638.

Berman, O., Krass, D., 2004. Facility location problems with stochastic demands and congestion, in: Drezner, Z., Hamacher, H. (Eds.), Facility Location: Applications and Theory, Springer. pp. 329–371.

Boffey, B., Galvao, R., Espejo, L., 2007. A review of congestion models in the location of facilities with immobile servers. European Journal of Operational Research 178, 643–662.

Castillo, I., A.Ingolfsson, Thaddues, S., 2009. Socially optimal location of facilities with fixed servers, stochastic demand, and congestion. Production and Operations Management 18, 721–736.

Elhedhli, S., 2005. Exact solution of class of nonlinear knapsack problems. Operations Research Letters 33, 615–624.

Elhedhli, S., 2006. Service system design with immobile servers, stochastic demand and congestion. Manufacturing and Service Operations Management 8, 92–97.

Gross, D., Harris, C.M., 1998. Fundamentals of Queueing Theory. 3 ed., John Wiley and Sons, New York.

Holmberg, K., Rönnqvist, M., Yuan, D., 1999. An exact algorithm for the capacitated facility location problems with single sourcing. European Journal of Operational Research 113, 544–559.

Huang, S., Batta, R., Nagi, R., 2005. Distribution network design: selection and sizing of congested connections. Naval Research Logistics 52, 701–712.

Marianov, V., Serra, D., 1998. Probabilistic, maximal covering location-allocation models for congested systems. Journal of Regional Science 38, 401–424.

Marianov, V., Serra, D., 2002. Location-allocation of multiple-server service centers with constrained queues or waiting times. Annals of Operations Research 111, 35–50.

Silva, F., Serra, D., 2008. Locating emergency services with different priorities: The priority queuing covering location problem. Journal of Operational Research Society 59, 1229–1238.

Vidyarthi, N., Elhedhli, S., Jewkes, E., 2009. Response time reduction in make-to-order and assemble-to-order supply chain design. IIE Transactions 41, 448–466.

Wang, Q., Batta, R., Rump, C.M., 2002. Algorithms for a facility location problem with stochastic customer demand and immobile servers. Annals of Operations Research 111, 17–34.

Zhang, Y., Berman, O., Macotte, P., Verter, V., 2010. A bilevel model for preventive healthcare facility network design with congestion. IIE Transactions 42, 865–880.

Zhang, Y., Berman, O., Verter, V., 2009. Incorporating congestion in preventive healthcare facility network design. European Journal of Operational Research 198, 922–935.

Zhang, Y., Berman, O., Verter, V., 2012. The impact of client choice on preventive healthcare facility network design. OR Spectrum 34, 349–370.