# An Expected Coverage Model with a Cutoff Priority Queue

Soovin Yoon

April 13, 2016

# 1 Introduction

# 2 Literature Review

## 2.1 Queuing Literature

### 2.1.1 Hypercube Literature

### 2.1.2 Cutoff Queuing Literature

## 2.2 Location Literatures with Priorities

# 3 Approximate Hypercube Model for Cutoff Priority Queue

In this section, we derive a Hypercube model approximation for cutoff priority queue. The Hypercube model can function as a standalone program for analysis of characteristics of cutoff priority queue. Also, the model iteratively interacts with MILP model described in the next section to assist decision-making for deployment and dispatch of ambulances under the uncertainty in their availability. The procedure introduced in this section extends the iterative procedure given by Larson(1975) as well as Jarvis(1985) for approximating equilibrium behavior and performance measures for the Hypercube model. As an extention, two additional features are implemented; calls are distinguished not only by their origins but also by their priorities, and low priority calls are "cut off" to reserve some servers for future high priority call arrivals.

A system with $s$ distinguishable servers and high(low) priority calls from the set of geographical locations $J$ is considered. High(low) priority calls for service arrive in a Poisson manner independent of busy status of servers, at a mean rate $\lambda^H(\lambda^L)$ per hour. The average service time for call from $j$ served by server $i$ is $\tau_{ij}$, which includes the travel time from the station to the scene, the time spent at the scene, transport time to the hospital and return travel time. Servers are distributed among the set of open stations $I^{open}$, with a
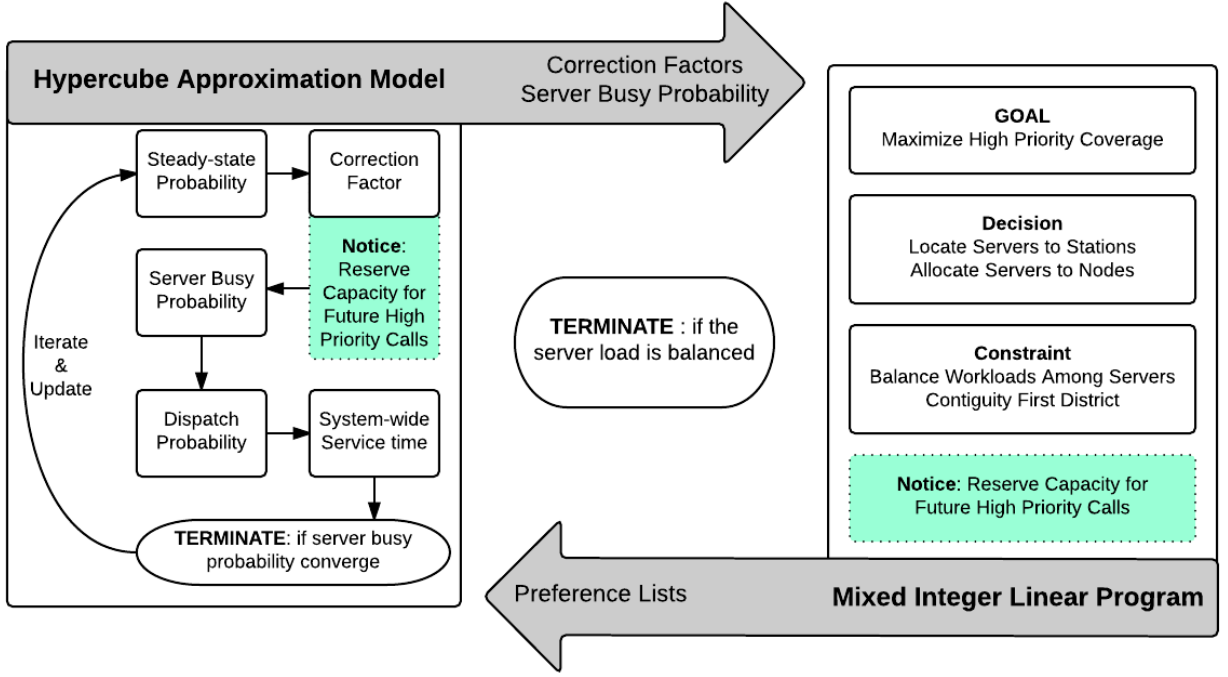
Figure 1: Iterative Procedure Flow Chart

server at a station. Exactly one server is assigned to a call immediately according to a fixed preference list, as long as there is any available server. In the preference list $\{b_{jk}^H\}(\{b_{jk}^L\})$,the $k$th preferred station by the high(low) priority call from node $j$ is given as $b_{jk}^H(b_{jk}^L) \in I^{open}$(for example, if $b_{12}^H = 3$ and $b_L 12 = 2$, then the second preferred server by the call from node 1 is 3 if the call is high priority and it is 2 if the call is low priority).

In order to implement the function of reserving resources on the system, a cutoff level of $s_1$ is given for admission of low priority customers. This means that low priority calls are placed in a queue or lost whenever $s_1$ or more servers are busy, instead of being served immediately. High priority calls are immediately served by the most preferred available server as long as there are free servers. They are either put in a queue or lost if all servers are busy. This implies that the server cutoff level for high priority customers is set to $s$ by convention(Larson 1986). The service is assumed to be non-preemptive.

We consider the two type of model that differs by the queue capacity. Under the assumption of $0-$line capacity system $(M/M/s/s)$, high priority calls arriving when all servers are busy or low priority calls arriving when there are more than or equal to $s_1$ busy servers are "lost to the system", as the queue has zero capacity. Lost calls are considered as being served by external resources. It is common in the EMS literature to assume the system with no queue, and it is also consistent with the real-world operation(Budge

2

| symbol | description |
|---|---|
| $j \in J$ | set of call nodes |
| $I \in I$ | set of potential stations |
| $i \in I^{open}$ | set of open stations, $I^{open} \subset I$ |
| $p \in \{H, L\}$ | set of call priorities |
| $s$ | total number of servers |
| $s_1$ | cutoff level for low priority calls |
| $\rho$ | server utilization |
| $r_i$ | busy probability of server at station $i$ |
| $r$ | system-wide server busy probability |
| $P_i$ | steady-state probability |
| $q_k$ | correction factor |
| $f_{jk}$ | probability of dispatching $k$th preferred server to call nodes $j$ |
| $b_{jk}$ | the station where the server $k$th preferred by call node $j$ is located |
| $\lambda_j^p$ | arrival rate of priority $p$ call from node $j$ |
| $\lambda^H(\lambda^L)$ | arrival rate of high(low) priority calls |
| $\lambda$ | system-wide call arrival rate |
| $\tau_{ij}$ | mean service time for calls from node $j$ served by a server from station $i$ |
| $\tau$ | system-wide mean service time |
| $\epsilon$ | server utilization change threshold |
| $\delta$ | server utilization imbalance threshold |

Table 1: Summary of Symbols

et al. 2009).

On the other hand, if the queue is assumed to have $\infty-$line capacity$(M/M/s/\infty)$, calls are put in a queue when they cannot be served immediately. They are served later in a first-come-first-served manner. As the queue is infinitely large, backlogged calls are never lost as they are eventually served in the long run. Here we show the derivation for both queue size assumptions.

Table 1 summarizes symbols used frequently. The hypercube queuing model requires deploying and dispatching policy as inputs, which can be generated either externally(for example, send-the-closest) or by solving the MILP in the next section. There are other parameters needed such as service times and arrival rates, and these can be collected from the data. The outputs of the model are correction factors $Q_k$ for $k = 1, \cdots, s$ and server busy probabilities $r_i$ for $i = 1, \cdots, s$. Then dispatch probability $f_{jk}$, the probability that a priority $p$ call from node $j$ is served by its $k$th preferred server, can be approximated as

$$f_{jk}^p = Q_{k-1}^p \Pi_{l=1}^{k-1} r_{b_{jl}^p} (1 - r_{b_{jk}^p}). \tag{1}$$

The correction factors $Q_k$ is needed to approximately correct for the unrealistic assumption of independent servers.

**STEP 0** Given:

$\lambda_j^p, \tau_{ij}$ from inputs, open stations $I^{open}$ and preference list $b_{jk}$

Initialize:

$$\tau^0 = \sum_{j \in J} \sum_{i \in I^{open}} \sum_{p \in \{H,L\}} \tau_{ij} \frac{\lambda_j^p}{\lambda^p} \tag{2}$$

$$\rho = \frac{\lambda \tau}{s} \tag{3}$$

**STEP 1** Update steady-state probabilities $P_0$ and $P_s$.

For $0-$line capaticy queue, use

$$P_i = P_0 \frac{s^i \rho^i}{i!}, \quad 0 \le i \le s_1$$
$$P_i = P_0 \frac{s^i \rho^i}{i!} \left(\frac{\lambda^H}{\lambda}\right)^{i-s1}, \quad s1+1 \le i \le s \tag{4}$$

where

$$P_0 = \left( \sum_{i=0}^{i=s1} \frac{s^i \rho^i}{i!} + \sum_{i=s1+1}^{s} \frac{s^i \rho^i}{i!} \left(\frac{\lambda^H}{\lambda}\right)^{i-s1} \right)^{-1}. \tag{5}$$

For $\infty-$line capacity queue, denoting $\rho_H = \lambda \tau P_H, \rho_L = \lambda \tau P_L, \rho = \rho_H + \rho_L$, use

$$P_n = P_0 \frac{\rho^s}{s!} \quad 0 \le n \le s_1 - 1 \tag{6}$$

$$= P_0 (\rho^{s_1} \rho_H^{n-s_1}/n!) s_1 [s_1 - \rho_L S_L(s_1, s)] \quad s_1 \le n \le s \tag{7}$$

where

$$P_0 = \left\{ \sum_{n=0}^{s_1-1} \rho^n/n! + (\rho^{s_1}/(s_1-1)!) S L(s_1,s)/[s_1 - \rho_L S_L(s_1,s)] \right\}^{-1}, \tag{8}$$

$$S_L(k,s) = \sum_{i=k}^{s} \rho_H^{i-k} k!/i!. \tag{9}$$

**STEP 2** Update correction factors $Q_k^H$ and $Q_k^L$ by

$$Q_k^H = \frac{\sum_{j=k}^{s-1} \frac{(s-k-1)!(s-j)j!}{(j-k)!s!} P_k}{w^k \rho^k \big(1 - (w\rho)\big)},$$

(10)

and

$$Q_k^L = \frac{\sum_{j=k}^{s_1-1} \frac{(s-k-1)!(s-j)j!}{(j-k)!s!} P_k}{w^k \rho^k \big(1 - w\rho\big)}.$$

(11)

For $0-$line capacity queue, use

$$w = (1 - P_s \frac{\lambda_H}{\lambda} - \sum_{i=s_1}^{s} P_i \frac{\lambda_L}{\lambda})$$

(12)

and $w = 1$ for $\infty-$line capacity queue.

**STEP 3** Update server busy probabilities $r_i$ as

$$r_i = \frac{V_i}{V_i + 1}.$$

(13)

where

$$V_i = \sum_{i \in I^{Open}} \sum_{k \in K} \sum_{p \in \{H,L\}} \lambda_j^p \tau_{b_{jk}^p j} q_k^p (\Pi_{l=1}^{k-1} r_{b_{jl}^p})$$

(14)

**STEP 4** Update system-wide service time $\tau$ as

$$\tau = \sum_{i \in I^{open}} \sum_{j \in J} \tau_{ij} \frac{\lambda_j^p}{\lambda} \frac{f_{j,a_{ij}^p}^p}{\sum_{k \in K} f_{jk}^p}$$

(15)

where $f_{jk}^H$ and $f_{jk}^L$ follows

$$f_{jk}^p = Q_{k-1}^p (1 - r_{b_{jk}^p})(\Pi_{l=1}^{k-1} r_{b_{jl}^p})$$

(16)

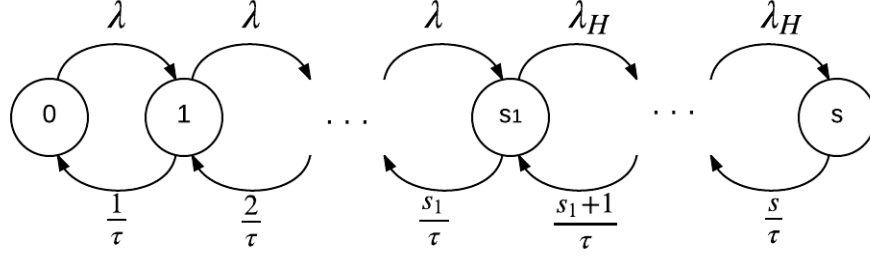**STEP 5** Update $\rho = \frac{\lambda \tau}{s}$ and $r = w\rho$.

Figure 2: Transition Diagram for 0−line Capacity Queue

Normalize $r_i$ by

$$r_i \leftarrow \frac{r_i}{\sum_i r_i/s} r \tag{17}$$

Check if max change in $r_i$ is below the predetermined threshold.

If yes, terminate with the return value $r$ and $Q_k^p$. If no, go back to step 1.

First, we need to derive closed-form expressions for steady-state probabilities in cutoff priority queue. For 0−line queue, Figure 2 shows the corresponding transition diagram.

Therefore, by solving the following balance equations, we gain the desired probability.

$$P_1 = \lambda \tau P_0 = \rho s P_0, \cdots, P_i = \frac{\rho^i s^i}{i!} P_0 \qquad i \leq s_1,$$

$$P_{s_1+1} = \lambda_H \frac{\tau}{s_1+1} P_{s_1} = \frac{\rho s}{s_1+1} \frac{\lambda_H}{\lambda} P_{s_1} = \frac{\rho^{s_1}}{s_1+1!} s^{s_1} \frac{\lambda_H}{\lambda} P_0,$$

$$\cdots, P_i = \frac{\rho^i s^i}{i!} (\frac{\lambda_H}{\lambda})^{i-s_1} P_0 \qquad i > s_1$$

In order to derive the steady-state probabilities formulation for the $\infty$−line queue, SCQL model in Taylor and Templeton(1980) was referred.

Next, correction factors for low priority in cutoff priority queue is derived. This follows Larson(1975)'s procedure to derive an expression for $P\{B_1 b_2 \cdots B_j F_{j+1}\}$ (p852) closely, punctuated by some changes due to cutoff assumption.

Let $B_j \equiv$ event that $j$th selected is busy (not available), $F_j \equiv B_j^c =$ event that $j$th server selected is free (or available), $S_k \equiv$ state of the queuing system indicating that exactly $k$ servers are busy, and

$P\{B_1 B_2 \cdots B_j F_{j+1}\} \equiv$ probability that the first free server is the $j + 1$st server selected.

Under the cutoff priority queue assumption, if the arriving call is low priority, we have

$$P\{B_1 B_2 \cdots B_j F_{j+1}\} = P\{B_1 b_2 \cdots B_j F_{j+1} | S_k, k < s_1\}$$

because we do not select a server to dispatch to a low priority call, if the number of busy server is above the cutoff.

Therefore we wish to derive an expression for $P\{B_1 b_2 \cdots B_j F_{j+1} | S_k, k < s_1\}$ that will motivate an approximation procedure for the model in which servers are not identical. (Note: In the hypercube model under a fixed-preference dispatching policy, the dispatcher always assigns the most preferred available server. Therefore the desired probability is the probability that the first $j$ preferred servers are busy and the $j + 1$st is free.)

By laws of conditional probability we have

$$P\{B_1 b_2 \cdots B_j F_{j+1} | S_k, k < s_1\} = \sum_{k=0}^{k=s_1-1} P\{B_1 B_2 \cdots B_j F_{j+1} | S_k\} P_k$$

Where

$$P\{B_1 B_2 \cdots B_j F_{j+1} | S_k\} = P\{F_{j+1} | B_1 B_2 \cdots B_j S_k\} P\{B_j | B_1 B_2 \cdots B_{j-1} S_k\} \cdots P\{B_1 | S_k\}.$$

Then we have

$$P\{B_i | B_1 B_2 \cdots B_{i-1} S_k\} = \frac{k - (i - 1)}{s - (i - 1)} \qquad i = 1, 2, \cdots, k + 1$$

$$P\{F_{j+1} | B_1 B_2 \cdots B_j S_k\} = \frac{s - k}{s - j} \qquad j = 0, 1, \cdots, k$$

Combining above results we have the desired probability

$$P\{B_1 b_2 \cdots B_j F_{j+1} | S_k, k < s_1\} = \sum_{k=j}^{s_1-1} \frac{k}{s} \frac{k-1}{s-1} \cdots \frac{k - (j-1)}{s - (j-1)} \frac{s-k}{s-j} P_k$$

$$= \sum_{k=j}^{s_1-1} \frac{(s - j - 1)!(s - k)k!}{(k - j)!s!} P_k$$

with $P_k$ derived in the note above, or,

$$P\{B_1 B_2 \cdots B_j F_{j+1} | S_k, k < s_1\} = Q_j^L r^j (1 - r),$$

where

$$Q_j^L = \frac{\sum_{k=j}^{s_1-1} \frac{(s-j-1)!(s-k)k!}{(k-j)!s!} P_j}{(1 - P_s \frac{\lambda_H}{\lambda} - \sum_{i=s_1}^{s} P_i \frac{\lambda_L}{\lambda})^j \rho^j \left(1 - (1 - P_s \frac{\lambda_H}{\lambda} - \sum_{i=s_1}^{s} P_i \frac{\lambda_L}{\lambda})\rho\right)}.$$

Lastly, the derivation of $Q_j^H$ is just analogous to above $s_1$ replaced by $s$.

## 4  Maximum Expected Coverage Model

In this section, a mixed integer linear program model that simultaneously locate servers at stations and dispatch servers to call nodes is introduced. The goal of the model is to maximize expected high priority coverage over the region. A spatial service system with $s$ servers, a set of potential stations $I$ and a set of demand nodes $J$ where calls arrive in a Poisson manner is considered. Exactly $s$ among $\|I\|$ stations are selected as open stations and one server is located for each open station. Also, each node $j$ is assigned a preference list $b_{jk}^H(b_{jk}^L)$, which is an ordered list of open stations for high(low) priority calls. A call is considered to be covered if a server is dispatched to the call immediately upon the arrival and it reaches the call in the predetermined time limit(e.g. 9 minutes).

Define the set of binary variable $y_i = 1(0), i \in I$ if a server is (not) located at station $i$. This variable set represents deployment decisions that defines the set of open stations. Also define the set of Bernoulli variables $x_{ijk}^p = 1(0), i \in I, j \in J, k = 1, \cdots, s$ if it is (not) the $k$th preferred station for a priority $p$ call from node $j$. This variable set captures the dispatch policy. Note that $\{x_{ijk}\}$ is basically a binary representation of the preference list $\{b_{jk}\}$ aforementioned.

The mixed integer linear program is formally stated below.

$$\max \sum_{i=1}^{s} \sum_{j \in J} \sum_{k=1}^{s} c_{ijk} x_{ijk}^H \tag{18}$$

subject to

$$\sum_{i \in I} y_i = s, \quad y_i \in \{0,1\} \tag{19}$$

$$\sum_{k=1}^{s} x_{ijk}^p = y_i \qquad \forall i,j,p \tag{20}$$

$$\sum_{i \in I} x_{ijk}^p = 1 \qquad \forall j,k,p \tag{21}$$

$$\sum_{p,j,k} \lambda_j^p q_k (1-r) r^{k-1} \tau_{ij} x_{ijk}^p \leq (r+\delta) y_i \tag{22}$$

8

$$\sum_{p,j,k} \lambda_j^p q_k (1-r) r^{k-1} \tau_{ij} x_{ijk}^p \geq (r-\delta) y_i \tag{23}$$

$$x_{ij'1}^p \geq x_{ij1}^p, \quad \forall p \in \{H, L\}, j \in J, i \in I, j' \in N_{ij} \tag{24}$$

where

$$c_{ijk} = q_{k-1}(1-r) r^{k-1} R_{ij} \frac{\lambda_j^H}{\lambda^H}. \tag{25}$$

The objective function in (18) represents expected coverage over high priority calls. The coefficient set $c_{ijk}$ used in the objective function is detailed in (25), and they represents explicit considerations for uncertainty in the ambulance unavailability and the travel time. Ambulance unavailability is captured by introducing $r$, the system-wide busy probability. Correction factor $q_k$, which is the output from previous Hypercube model, yields approximation to actual queuing probabilities that compensate unrealistic independence assumption between the servers. Non-deterministic travel time is also implemented by adopting a parameter $R_{ij}$, which represents the probability that an ambulance located at station $i$ may reach the node $j$ within the time limit. $R_{ij}$ is naturally a decreasing function of the distance between the station $i$ and node $j$. Putting the parameters together, $q_{k-1}(1-r) r^{k-1} R_{ij}$ captures the coverage of the server at station $i$ over high priority calls from node $j$ when $i$ is the $k$th preferred server by high priority call at $j$. And then this value is weighted by its relative call volume proportion $\lambda_j^H/\lambda^H$ to shape the coefficient $c_{ijk}$, so that the objective function can represents expected coverage over the entire region.

A station is considered to be opened if a server is located, and there are exactly $s$ open stations which is guaranteed by (19). (20) ensures that a server from station $i$ is dispatched only when the station $i$ is open. Also, (21) enforces all call nodes to have an ordered, non-overlapping preference list.

(22) and (23) are balancing constraints that sets upper and lower limits for server utilization. The server utilization is kept close to average server utilization within the threshold of $\delta$ for all servers. This constraint set has two implications. First, the program uses common busy probability $r$ instead of $r_i, i \in I^{open}$ implying the implicit assumption that the busy probability is the same for all ambulances. This assumption is reasonable in an endogenous sense, as the load balancing is approximately guaranteed by (22) and (23) to keep the server utilization equal. Also, this constraints are meaningful in the sense of equity from the service provider's perspective. They ensures that paramedics and emergency medical technicians who staff each server avoid exposure to distress by excessive workload but also have opportunities to practice their skills(McLay and Mayorga 2013).

(24) refers to contiguity constraints, which enforces the program to assign the first priority in a geographically reasonable sense. When the first preferred station for high priority calls from node $j$ is station $i$, and there is a node $j'$ is a neighborhood of $j$ and closer to $i$ than $j$, then constraint set (24) enforces $i$ to be the

first preferred station for a high priority call from node $j$.

Finally, note that the cutoff priority queue scheme is brought to the MILP model through the value of $q_k^L$. $q_k^L$, which is an output from the Hypercube model, is set to zero for all $k < s_1$. In this program, when $i$ is the $k$th preferred server by priority $p$ call from node $j$, the probability of dispatching a server $i$ to priority $p$ call from node $j$ is computed as $q_k^p(1-r)r^{k-1}$. Therefore, the probability of dispatching $k$th preferred servers to low priority calls are always zero.

Once the MILP problem is solved, the preference list $b_{jk}^p$ is created from the optimal value of $x_{ijk}^p$ accordingly by Algorithm 1.

**for** $j \in J$ **do**
    **for** $i \in I^{open}$ **do**
        $k \leftarrow 1$ **while** $k \leq s$ **do**
            **if** $x_{ijk}^H = 1$ **then**
                | $b_{jk}^H \leftarrow i$;
            **end**
            **if** $x_{ijk}^L = 1$ **then**
                | $b_{jk}^L \leftarrow i$;
            **end**
            $k \leftarrow k+1$
        **end**
    **end**
**end**

**Algorithm 1:** Pseudo-code for Generating Preference Lists

The system-wide server busy probability $r$ is used as a parameter in objective coefficients as well as load balancing constraints. However, this is an endogenous value in that the decision of deployment and dispatching itself affects the busy probability of servers. Nevertheless, in pursuit of problem tractability, we want to avoid expressing $r$ as a function of decision variable $y_i$ and $x_{ijk}^p$, because doing so would make the mixed integer program non-linear.

In order to overcome this complication, an iterative scheme is introduced. As aforementioned in the previous section, the Hypercube model is utilized to generate the correction factors $q_k$ and server busy probabilities $r$ used as inputs for MILP model. Then the solutions from the MILP model become the new inputs for the Hypercube model to update the value of correction factors and server busy probabilities. Again, those outputs are utilized as inputs for the MILP model. This iterative procedure continues until outputs converges so that the change in server busy probability drops under the threshold.

Table (2) details the steps in iterative Algorithm that comprehensively describes the iterative use of

**STEP 0 Initialize**:

Generate the initial preference lists $b_{jk}^p$ from the *send-the-closest-available-server* rule.

Set correction factors $q_k^p = 1$ for all $k = 1, \cdots, s$, server busy probability $r = \rho = \dfrac{\lambda \tau}{s}$, where system-wide mean service time $\tau$ set as $\tau = \lambda^{-1} \sum_i \sum_j \sum_p \lambda_j^p \tau_{ij}$.

**STEP 1** Solve the MILP model with inputs $r, r_i, q_k^p$ to update preference lists $b_{jk}^p$ and set of open stations $I^{open}$.

**STEP 2** Solve the Hypercube model with inputs $b_{jk}^p$ and $I^{open}$ to update server busy probabilities $r, r_i$ and correction factors $q_k^p$.

**STEP 3** Check termination criteria. Stop if (1) max change in $r_i$ is less than the imbalance threshold or if (2) MILP was infeasible.

If both termination criteria are not satisfied, go back to **STEP 1**.

Table 2: Iterative MILP-Hypercube Algorithm

Hypercube model and MILP model.

# 5 Computational Results

This section reports analytical and simulation results that illustrates effectiveness of introducing the cutoff priority queue scheme for urban emergency medical service system. A loss system is considered mainly, but computational results for infinite-line queue are also introduced briefly for comparison.

## 5.1 Computational Setting

The model is investigated by using real-world data set from Hanover County, Virginia. The data describes a geographical region with 16 potential stations, 270 aggregated demand nodes with 139 2×2 mile cells and 131 1×1 mile cells. Call arrival rate during the weekdays 12PM to 6PM was selected for demand mean. The demand volume over the region is visualized in Figure 3 as a colored cells on a grid that darker color represent higher call volume as it can be seen from the color bar on the right. The call volume over the area is not uniform, as it is concentrated in the middle and south-east area of the county while the call is relatively scarce on the west side. The location of 16 potential stations are also marked in the map.

In pursuit of investigating the effect of reserving capacities in the large fleet system, a scaled-up version of Hanover County dataset was also used. The number of server was increased to 16 from 5, and the demand data was scaled up by a factor of 5.9. Note that the demand wasn't scaled up by the ratio of fleet size(3.2=16/5), but instead the scale parameter was selected so that the system show the similar coverage
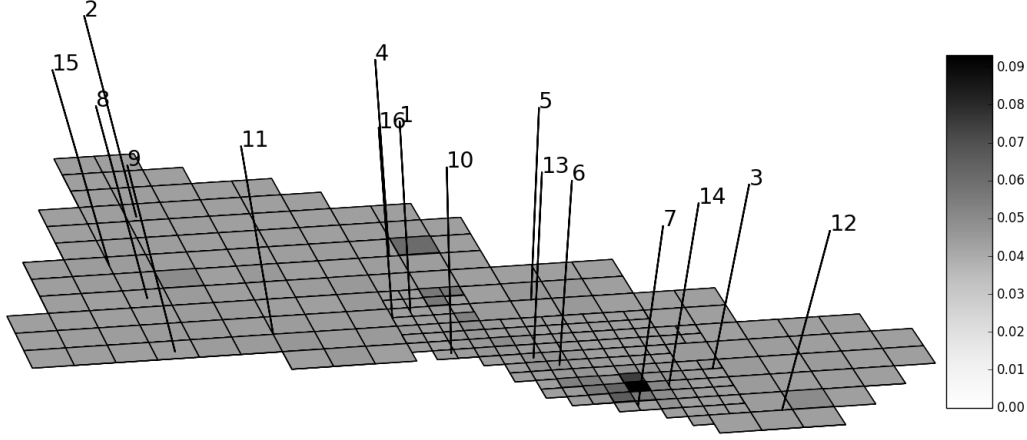
Figure 3: Distribution of Call Demand and Station Location on the Hanover County Map

performance to 5 server case when no server is reserved.

The iterative model was implemented in Python 3.4, and Gurobi 6.5.0 is used as an solver for the Mixed Integer Linear Program. The running time for the iterative MILP-Hypercube procedure using Hanover County data were 27 seconds in average including the time spend for running Hypercube models less than 0.17 second for all cases. The runtime for scaled data with large fleet was The load imbalance tolerance threshold was set to 2 percent, and server busy probability imbalance tolerance was set to 1 percent.

Throughout the section, analytical results were compared to the result of discrete-event simulation. The purpose of evaluating the result by simulation is to confirm the validity of cutoff model as well as iterative MILP-Hypercube procedure. The simulation used the same demand distribution and geographical data as the analytical model, and the resulted deployment and dispatch decisions from the analytical model were also used as inputs for the simulation to set the policy. The program was implemented in Python 3.4, ran until 100,000 call arrives and repeated 30 times to generate summary statistics for dispatch probabilities, server busy probabilities, steady-state probabilities and coverage over the high priority calls. The computational time for the simulation were $-,-$ and $-$ seconds for cutoff values $s1 = 3, 4, 5$ each. Note that the runtime of the simulation is significantly longer than the runtime for Hypercube model that serves the similar function, which is to evaluate the system given the preference lists. In pursuit of precisely getting statistics from rare events with less than 1 percent($P_s$ from the event that all servers are busy, $f^p(jk)$ for $k = 5$, the simulation should be run for a long time, which highlights the another benefit of developing an analytical model like the Hypercube model.
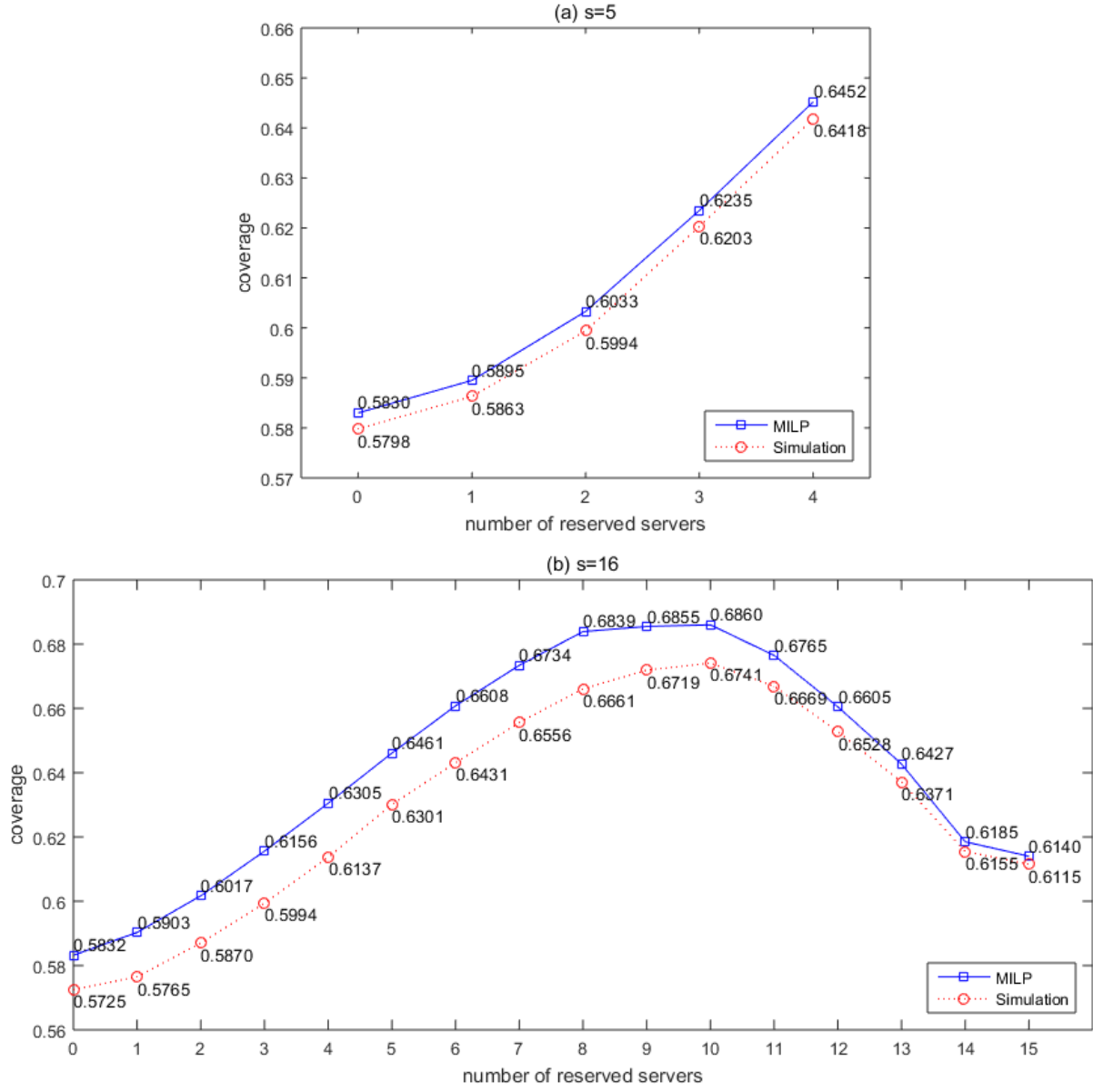
Figure 4: Coverage over High Priority Calls Using Different Cutoff Numbers in (a) Hanover County dataset (b) Scaled Hanover County dataset with Large Fleet

## 5.2 Coverage Improvement

The resulted coverages over the high priority calls are displayed in figure 4. The left plot (a) in figure 4 show the resulted expected coverage over high priority calls from MILP model and Simulation with 5 servers. The coverage improvement is monotonic as the system reserve more servers by setting the cutoff number $s_1$ lower.

The right plot (b) was drawn over a scaled-up version of Hanover County dataset. The effect of capacity reservation on high-priority call coverage is more clear in this large fleet system. Specifically, always reserving 8 half of the fleet for high priority calls, the coverage improvement is as large as – percent. The marginal improvement was largest when the number of reserved server was increased from 3 to 4, that the decrement of cutoff number $s_1$ just by 1 resulted the coverage improvement of 1.49 percent. Considering the fact that the maximum achievable coverage of this system is only –(which is gained by assuming that every call is served immediately by the server that is located at the most reachable station), the improvement from – to – by reserving only one server is significant.

Another noticeable observation is that there exist a right level of capacity reservation. The coverage improvement appears to be not monotonic in number of reserved servers. (better use a expression like threshold/upperbound/cutoff?)

(TODO) Add a plot for simulation result that keep using the same preference lists, but show how imbalanced their server utilization is. That will explain why coverage is not improving over a certain reserving number(or under a certain cutoff number)

Lastly, in a model validation context, the analytical results were compared to the simulation results. The dotted lines in figure 4 (a) and (b) shows the simulation result that corresponds to the analytical result. The coverage levels for all cutoff numbers were very close for both of analytical MILP model and the simulation model, with differences in coverage result less than – percent in the original data and less than – percent in the large fleet system, for all cutoff number $s_1$ given.

(Q) How about reordering subsections: showing coverage first and hypercube statistics later? Because we have to talk about the trade-off between loss in p(serve low calls) and gain(coverage improvement)

The resulted effect on the coverage is a collaboration of two factors. First, by reserving more servers for high priority call arrivals, the dispatch probability of sending highly preferred server increases, which results in the improvement in the coverage. Also, with different correction factor and server busy probability given MILP program produces different deployment and dispatch decisions that is optimized for given cutoff number. We will examine each factor in detail in the next subsection.

| | Hypercube | | | | | Simulation | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | $s_1=1$ | $s_1=2$ | $s_1=3$ | $s_1=4$ | $s_1=5$ | $s_1=1$ | $s1=2$ | $s_1=3$ | $s_1=4$ | $s_1=5$ |
| k=1 | 0.7918 | 0.7371 | 0.6980 | 0.6771 | 0.6698 | 0.7887 | 0.7312 | 0.6950 | 0.6721 | 0.6664 |
| k=2 | 0.1672 | 0.1949 | 0.2033 | 0.2023 | 0.2000 | 0.1628 | 0.1871 | 0.1987 | 0.2011 | 0.1979 |
| k=3 | 0.0305 | 0.0507 | 0.0656 | 0.0699 | 0.0690 | 0.0356 | 0.0545 | 0.0689 | 0.0712 | 0.0707 |
| k=4 | 0.0069 | 0.0115 | 0.0220 | 0.0279 | 0.0276 | 0.0087 | 0.0143 | 0.0241 | 0.0311 | 0.0301 |
| k=5 | 0.0020 | 0.0033 | 0.0062 | 0.0128 | 0.0126 | 0.0028 | 0.0043 | 0.0083 | 0.0142 | 0.0135 |
| $P(\text{lost})$ | 0.0015 | 0.0025 | 0.0049 | 0.0100 | 0.0209 | 0.0015 | 0.0025 | 0.0050 | 0.0102 | 0.0214 |

Table 3: High Priority Dispatch Probabilities for Different Cutoff Numbers

| | Hypercube | | | | | Simulation | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | $s_1=1$ | $s_1=2$ | $s_1=3$ | $s_1=4$ | $s_1=5$ | $s_1=1$ | $s1=2$ | $s_1=3$ | $s_1=4$ | $s_1=5$ |
| k=1 | 0.2813 | 0.5160 | 0.6296 | 0.6646 | 0.6698 | 0.2833 | 0.5157 | 0.6300 | 0.6605 | 0.6701 |
| k=2 | 0.0000 | 0.0741 | 0.1503 | 0.1897 | 0.2000 | 0.0000 | 0.0748 | 0.1489 | 0.1941 | 0.1983 |
| k=3 | 0.0000 | 0.0000 | 0.0280 | 0.0572 | 0.0690 | 0.0000 | 0.0000 | 0.0285 | 0.0568 | 0.0698 |
| k=4 | 0.0000 | 0.0000 | 0.0000 | 0.0152 | 0.0276 | 0.0000 | 0.0000 | 0.0000 | 0.0152 | 0.0276 |
| k=5 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0126 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0126 |
| $P(\text{lost})$ | 0.7185 | 0.4093 | 0.1915 | 0.0729 | 0.0209 | 0.7158 | 0.4094 | 0.1923 | 0.0737 | 0.0214 |

Table 4: Low Priority Dispatch Probabilities for Different Cutoff Numbers

| | Hypercube | | | | | Simulation | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | $s_1=1$ | $s1=2$ | $s_1=3$ | $s_1=4$ | $s_1=5$ | $s_1=1$ | $s1=2$ | $s_1=3$ | $s_1=4$ | $s_1=5$ |
| $r$ | 0.2082 | 0.2629 | 0.3020 | 0.3229 | 0.3302 | 0.2060 | 0.2618 | 0.3014 | 0.3224 | 0.3300 |

Table 5: Server Utilization

## 5.3 Dispatch Probability

Dispatch probabilities for loss systems are shown in Table 5.3 and 3. Dispatch probabilities for higher $k$ $(k = 1, 2)$ for high priority calls increased as the system reserve more servers(lower cutoff number $s_1$). Dispatch probabilities for $k > s_1$ for low priority calls are all zero, which is the direct result of cutoff priority queue discipline.

Finally, considering the fact the Hypercube model is an approximate model, the discrepancy between the analytical Hypercube model results and simulation results are relatively small enough to conclude that the Hypercube model is valid.

## 5.4 Infinite-line Queue Systems

(TODO) Brief results under the assumption of infty queue, to show there's no noticeable improvement under this scheme, so that we can validate loss system assumption. Let's simply display a dispatch probability table and coverage graph here which was drawn in the exact same way as previous sections, so that it clearly shows that with infinity queue reserving servers does not mean that much. Or what can it mean if it does.