

Authors are encouraged to submit new papers to INFORMS journals by means of a style file template, which includes the journal title. However, use of a template does not certify that the paper has been accepted for publication in the named journal. INFORMS journal templates are for the exclusive purpose of submitting to an INFORMS journal and should not be used to distribute the papers in print or online or to submit the papers to another publication.

A Bound on the Performance of Optimal Ambulance Redeployment Policies in Loss Systems

Kenneth Chong, Shane Henderson, Mark Lewis, Huseyin Topaloglu
School of Operations Research and Information Engineering, Cornell University, Ithaca, New York
kcc66, sgh9, mel47, ht88@cornell.edu

Ambulance redeployment refers to the practice of strategically relocating ambulances in real time to better respond to future call arrivals. Optimal redeployment policies can be found, in theory, via stochastic dynamic programming, but this approach succumbs to the curse of dimensionality. We instead consider the problem of finding an upper bound on the performance of optimal redeployment policies, which can be used, for instance, as a benchmark against which heuristic policies can be compared. Such a bound has been developed for systems in which calls can queue when no ambulances are available. We adapt this bound to loss systems, corresponding to the fairly common situation in which arriving calls are redirected to an alternative service if they cannot be immediately served. This adaptation is nontrivial, and involves the introduction of new ideas, including an integer program whose feasible region contains the set of decisions made by any member of a fairly general class of redeployment policies. We prove the validity of our upper bound, and study its performance through computational experiments.

Key words: ambulance redeployment; system status management; dynamic programming; upper bounds; integer programming; loss systems; curse of dimensionality

1. Introduction

Emergency Medical Service (EMS) providers are tasked with responding to calls for assistance, and must often do so while contending with factors such as rising demand, increasing medical costs, and traffic congestion (Henderson 2011). To attain a high level of service in such an environment, EMS providers can make a number of operational decisions to improve performance. Methods

for doing so include changing how ambulances are assigned to base locations, or how ambulances are dispatched to incoming calls. More recently, EMS providers have also adopted the practice of *ambulance redeployment*, which entails strategically repositioning idle ambulances (either from home bases or following service completions) to improve the system's responsiveness to future calls (Stout 1989). Although the benefit associated with implementing redeployment policies tends to be modest (Alanis et al. 2013, Maxwell et al. 2014), it is practically significant, and the practice has been adopted by a large fraction of EMS systems in North America (Williams 2009).

Although the problem of finding an optimal redeployment policy can be formulated as a stochastic dynamic program, and has been done, for instance, in Berman (1981a,b,c) and McLay and Mayorga (2013), such a model would need to include the location and status of every ambulance in the state space. This approach is impractical in large-scale systems, as it succumbs to the curse of dimensionality, and would likely yield a policy that is too unwieldy to implement. As a result, a number of methods for obtaining good policies have been proposed, and techniques for doing so include using approximate dynamic programming (Maxwell et al. 2010, Schmid 2012), solving integer programs in real time (Gendreau et al. 2001, Nair and Miller-Hooks 2009, Naoum-Sawaya and Elhedhli 2013, Zhang 2012), developing heuristics (Andersson and Vaerband 2015, Jagtenberg et al. 2015, van Barneveld et al. 2015, 2016), and building compliance tables, which specify where ambulances should be positioned, given coarse information about the state of the system (Alanis et al. 2013, Gendreau et al. 2006, Sudtachat et al. 2016, van Barneveld 2016).

However, an important question to ask is whether any of these heuristic policies can be shown to be near-optimal. To address this concern, methods have been proposed to construct bounds on the performance of a fairly general class of redeployment policies; see Yue et al. (2012) and Maxwell et al. (2014). These bounds are typically not tight, but they may still be informative. For instance, they allow a municipality to determine if performance targets are achievable by a proposed ambulance schedule, which was the primary motivation for Maxwell et al. (2014). Alternatively, bounds provide an indication of the increase in performance that can be realized by implementing

redeployment policies. Because these policies remain controversial, due to the heavier burden they place on emergency medical staff (Bledsoe 2003), bounds may help in determining whether such an operational change is worthwhile.

Computational experiments suggest that the bound in Maxwell et al. (2014) is tighter than that in Yue et al. (2012), but the former bound assumes that arriving calls can queue when all ambulances are busy. Many EMS systems make use of external resources (such as police and fire vehicles, or ambulances from other municipalities) during “red alert” situations. This is the case in some large-scale systems (myFoxDetroit 2011, Pedersen 2012, Sinnema 2010, 2012), as well as many small- to medium-scale systems that pool resources with neighboring EMS providers. When this is the case, it may be more reasonable to model the EMS as a *loss system*. This is a subtle difference, but presents a significant technical challenge, as the bound in Maxwell et al. (2014) is not valid for such systems.

In this paper, we propose a new upper bound that is provably valid in loss systems. Although we draw upon ideas in Maxwell et al. (2014) to do so, significant new ideas are required. In contrast to Maxwell et al. (2014), and like in Yue et al. (2012), we obtain an upper bound via optimization along each sample path, by constructing an integer program for which the decisions made by almost any redeployment policy correspond to a feasible solution. Since we model the EMS as a loss system, we cannot directly compare our bound to that by Maxwell et al. (2014), but computational experiments suggest that our upper bound is tighter than one originating from a perfect information relaxation of the problem in the style of Yue et al. (2012).

The remainder of the paper is organized as follows. In Section 2, we explicitly formulate our model of the ambulance redeployment problem. In Section 3, we review the concepts from Maxwell et al. (2014) that we use in developing our upper bound, and demonstrate by counterexample that their bound is not valid for loss systems. We introduce our modification to their upper bound in Section 4, and prove its validity. Finally, we numerically evaluate our upper bound in Section 5, and conclude in Section 6.

2. Problem Formulation

Consider an EMS system operating A ambulances, whose service area is represented by a graph $G = (V, E)$, where V is a set of demand nodes, and E is a set of edges. Ambulances can be deployed to a set of base locations $B \subseteq V$. Let t_{ij} denote the travel time between demand node i and base j ; we assume this quantity is deterministic and time-stationary. Calls arrive according to an arbitrary stochastic process that we assume is independent of the system state, as well as of any decisions made by the EMS provider. We further assume that the locations of these calls are i.i.d. random variables; let d_i denote the probability that an arriving call originates from node $i \in V$.

Ambulances must immediately be dispatched to arriving calls whenever one is available, but we place no constraints on which ambulance must be sent. We say that a call receives a *timely response* if an ambulance reaches the call's location within a response threshold τ_r (typically, 9 minutes in practical applications). The response interval typically includes the time required for the dispatcher to triage and assign an ambulance to the call, and for the crew in the responding ambulance to prepare for travel to the scene, which we refer to as the *chute time*. If a call receives a response (regardless of whether or not the response is timely), the ambulance travels to the call's location, and spends a random amount of time on scene treating the patient; let F denote the distribution function of this random variable. We assume, for convenience, that calls do not require transport to a hospital following treatment on scene (and so service completes at the call's location), but this assumption can be easily relaxed, as we briefly explain later. Calls arriving when all ambulances are busy leave the system.

At any time, the EMS provider may redeploy idle ambulances, so as to increase the likelihood of timely responses to future call arrivals, and may do so in an arbitrary fashion. A redeployment move entails relocating an ambulance from a node i to a base j ; we need not specify the process by which this move occurs, provided that its location at any point is a node in the graph. Specifically, our upper bound is valid if we assume that the ambulance idles at node i for t_{ij} time units, then immediately relocates to base j , or makes a sequence of smaller "jumps" within the graph, or any

variant thereof, but we do not allow for “continuous” motion within the graph. The only restriction we place on the redeployment policies we consider is that they do not anticipate— that is, they do not use information about future call arrivals.

We evaluate the performance of a policy via the expected number of calls receiving timely responses over a finite horizon $[0, T]$, which, by the strong law of large numbers, is roughly equivalent (for large enough T) to the objective of maximizing the long-run average proportion of timely responses to calls. For convenience, we will sometimes refer to the former quantity as the total reward collected by the system.

3. Preliminaries

We begin by reviewing the tools from Maxwell et al. (2014) that we use in developing our bound.

3.1. A Bounding Reward Function

Let $v : \{0, 1, \dots, A\} \rightarrow [0, 1]$ be a nondecreasing function for which $v(a)$ is an upper bound on the probability that a call arriving when a ambulances are available receives a timely response. In practice, this probability depends both upon the locations of the a ambulances, and upon the call’s arrival time, but $v(a)$ is a valid upper bound regardless of these factors. We let $v(0) = 0$ in our model, as calls leave the system when no ambulances are available. (This is in contrast to Maxwell et al. (2014), who let $v(0) = v(1)$, as a call arriving when $a = 0$ may still receive a timely response if a busy ambulance is about to become available.)

To see how this function can be used to construct an upper bound, consider an alternate EMS system in which the decision maker collects a reward $v(a)$ when a call arrives while a ambulances are available (instead of a reward of either 0 or 1, depending on whether the response is timely). Maxwell et al. (2014) show that each redeployment policy π , in expectation, performs at least as well in the alternate system as in reality. More formally, let R^π and V^π be the reward collected by policy π in the original and alternate systems, respectively.

LEMMA 1. *For every policy π , $E[R^\pi] \leq E[V^\pi]$.*

See Maxwell et al. (2014) for a proof. They also demonstrate that for a given a , $v(a)$ can be obtained by solving an instance of the Maximal Covering Location Problem (MCLP), which is due to Church and ReVelle (1974). This entails solving an integer program which locates a ambulances on a graph so as to maximize the proportion of calls “covered” by an ambulance, where a demand node is considered to be covered if the closest ambulance can provide a timely response.

3.2. Bounding Service Time Distributions

Next, let $\{G_a : 0 \leq a \leq A\}$ be a collection of distribution functions, where $G_a : [0, \infty) \rightarrow [0, 1]$ is a stochastic lower bound on the service time distribution associated with a call, provided that a ambulances are available when it arrives (regardless of their locations). Here, we define service time to be the sum of the chute time, the travel time to the call’s location, and the on-scene treatment time. (In our case, G_0 is a Dirac delta function centered at 0, but we define $G_0 := G_1$ to be consistent with Maxwell et al. (2014).) Assume the distributions are stochastically nonincreasing in a — that is, for all $x \geq 0$, we have that $G_0(x) \leq G_1(x) \leq \dots \leq G_A(x)$. Maxwell et al. (2014) observe that given x and a , $G_a(x)$ can be computed by solving a p -median problem, which locates a ambulances within a graph so as to maximize the probability that the service time associated with the next call is at most x time units. They formulate this problem as an integer linear program.

Remark: Maxwell et al. (2014) show the p -median problem can be formulated so that the distributions $\{G_a : 0 \leq a \leq A\}$ are stochastic lower bounds on the service time distribution even when patients may require transport to a hospital. Let $H \subseteq V$ be the set of hospitals in the system, q_{ih} be the probability that a call arriving to node i is transported to hospital h (we may have $\sum_{h \in H} q_{ih} < 1$ for some i), and F_h be the distribution of the patient transfer time at hospital h . Defining service time to also include any time needed to transport the patient to the hospital, as well as to complete the transfer, Maxwell et al. (2014) numerically convolve F with the F_h to obtain the service time distribution associated with an arrival to node i . These distributions are used to obtain inputs for the p -median problem.

3.3. A Counterexample

To obtain an upper bound that is valid for any redeployment policy π , Maxwell et al. (2014) construct a queueing system, which they call their *bounding system*, with A servers (ambulances) in which jobs (calls) arrive according to the same stochastic process as in the original EMS system. Rewards and service times are state-dependent, in that a call arriving when a ambulances are idle generates reward $v(a)$, and spends a random amount of time in service that is governed by the distribution function G_a . Calls arriving when all ambulances are busy are placed in an infinite buffer, and served in a first-in, first-out fashion.

They demonstrate using a coupling argument that the reward collected by the bounding system, in expectation, is at least as large as the performance attained by an optimal policy in the original EMS system. Suppose that N denotes the number of calls arriving in the interval $[0, T]$, that they arrive at times T_1, \dots, T_N , and that A_k^π represents the number of ambulances available at time T_k under policy π . Let \tilde{A}_k be analogous to A_k^π for the bounding system. The validity of their upper bound follows from Lemma 1, and from showing that $\tilde{A}_k \geq A_k^\pi$ holds pathwise for all π and k .

In loss systems, the above relation may not hold along every sample path. There may be times at which the bounding system can respond to a call, but the original EMS system may be forced to turn away the call, leading to work that only the bounding system is forced to process. Thus, there may exist sample paths ω on which $\tilde{A}_k(\omega) < A_k^\pi(\omega)$ for some k . This suggests there may be loss systems for which the bound is not valid, and we demonstrate this by counterexample:

EXAMPLE 1. Consider a graph with two nodes, connected by an edge which takes 1 minute to traverse. Let $d_1 = d_2 = 0.5$, implying that a call is equally likely to arrive at either node. In this system, $A = 2$, $T = 60$, and at the start of the horizon, one ambulance is stationed at each node. The arrival process is deterministic, with calls arriving at 8, 16, 24, 29, 38, and 40 minutes after the start of the horizon. The response time threshold is set to 0 minutes, and so timely responses can only be provided from ambulances situated at the call's location. On-scene treatment times are also deterministic, and last 10 minutes.

Suppose the EMS provider dispatches the closest available ambulance to incoming calls, and that never redeploys ambulances following service completions. Then a direct calculation shows that this policy, in expectation, provides 3.25 timely responses. Computing the bound from Maxwell et al. (2014), we find that $v(0) = v(1) = 0.5$ and $v(2) = 1$, that G_0 and G_1 correspond to random variables that take on the values 10 and 11 with probability 0.5, and that G_2 is a Dirac delta function centered at 10. Using simulation or analytical methods, we find their method yields an upper bound of 3. Since we have found a policy that attains a strictly higher objective value, the upper bound is not valid.

4. A Modified Upper Bound

4.1. Construction

To develop an upper bound that is valid for loss systems, we construct a random variable Z such that when it is jointly defined with V^π on the same probability space, Z dominates V^π pathwise under any policy π . We begin by specifying the sample space Ω associated with our model. Random quantities in our model include N , the call arrival times T_1, \dots, T_N , the locations L_1, \dots, L_N to which these calls arrive, and the corresponding on-scene treatment times. We characterize the randomness associated with the latter by tagging the k^{th} arriving call with a Uniform(0, 1) random variable U_k , from which the corresponding on-scene treatment time can be obtained by computing $F^{-1}(U_k)$. While this characterization is somewhat indirect, it facilitates our construction of Z . A sample path thus consists of realizations of N , the total number of arriving calls, and of T_k , L_k , and U_k for each call k .

As in Maxwell et al. (2014), we also construct a bounding system with A ambulances, and in which calls arrive according to the same stochastic process. Rewards and service times are state-dependent in the same way as specified in Section 3.3. Contrary to Maxwell et al. (2014), calls arriving when all ambulances are busy are lost, the full sample path is revealed to the decision-maker at the beginning of the horizon, and calls can be rejected from the system, even when ambulances are available. The latter assumption may appear counterintuitive, but it may be preferable to

respond to fewer calls, and to collect the larger rewards associated with admitting these calls in “higher” system states. The bounding decision-maker seeks a “policy”—that is, a subset of calls to admit into the system—that maximizes total reward. We define $Z(\omega)$ to be the reward collected by an optimal policy in the bounding system under sample path ω .

Given such an ω , we can compute $Z(\omega)$ by solving an integer program. Let v and $\{G_a : 0 \leq a \leq A\}$ be as in Section 3. Suppressing ω for clarity, we let x_{ka} be a binary variable that takes on the value 1 if call k is served when a ambulances are available, and y_k denote the number of available ambulances when the k^{th} call arrives. Finally, let Q_k be a set of pairs (ℓ, a) such that if the ℓ^{th} call is admitted when a ambulances are available, service of the call completes after time T_{k-1} , but before time T_k (again, with ω suppressed for clarity). More formally, $(\ell, a) \in Q_k$ iff $T_{k-1} < T_\ell + G_a^{-1}(U_k) \leq T_k$. Our formulation is as follows:

$$Z(\omega) := \max \sum_{k=1}^N \sum_{a=1}^A v(a) x_{ka} \quad (\text{IP}(\omega))$$

$$\text{s.t.} \quad \sum_{a=1}^A x_{ka} \leq 1 \quad \forall k \in \{1, \dots, N\} \quad (1)$$

$$\sum_{a=1}^A a x_{ka} \leq y_k \quad \forall k \in \{1, \dots, N\} \quad (2)$$

$$y_k = y_{k-1} - \sum_{a=1}^A x_{k-1,a} + \sum_{(\ell, a) \in Q_k} x_{\ell,a} \quad \forall k \in \{2, \dots, N\} \quad (3)$$

$$y_1 = A \quad (4)$$

$$x_{tk} \in \{0, 1\}, \quad y_k \in \{0, 1, \dots, A\}$$

Constraints (1) enforce, for each call k , that $x_{ka} = 1$ for at most one choice of a . Constraints (2) serve two purposes: ensuring that no calls are admitted when all ambulances are busy, and setting $x_{k,y_k} = 1$ whenever the k^{th} arriving call is admitted. The latter occurs naturally because the objective coefficients $v(a)$ are nondecreasing in a , and so we would not have $x_{ka} = 1$ for some $a < y_k$. Constraints (3) preserve “flow balance”: the number of idle ambulances when the k^{th} call arrives matches the number of ambulances available at time T_{k-1} , minus one if an ambulance was

dispatched to call $k - 1$, plus the number of ambulances that become idle between time T_{k-1} and time T_k . Finally, constraint (4) enforces that all ambulances are initially idle.

Remark: The integer program (IP(ω)) can be viewed as an adaptation of that by Yue et al. (2012), who also consider a relaxed problem in which the decision-maker has perfect information. Although our model, in contrast to theirs, does not consider the locations of ambulances in the bounding system, knowing the number of available ambulances at a given point in time suffices to compute our upper bound.

4.2. Validity for Loss Systems

To verify that Z can be used to upper bound the performance of any redeployment policy in loss systems, we prove the following:

LEMMA 2. *For every sample path ω and policy π , $V^\pi(\omega) \leq Z(\omega)$.*

Proof: Fix π and ω . Let $S^\pi(\omega) \subseteq \{1, 2, \dots, N(\omega)\}$ denote the set of calls to which π dispatches an ambulance on this path (in the original system). We have that

$$V^\pi(\omega) = \sum_{k \in S^\pi(\omega)} v(A_k^\pi(\omega)).$$

We use $S^\pi(\omega)$ to construct a feasible (but not necessarily optimal) solution (\bar{x}, \bar{y}) to (IP(ω)) that attains an objective value of at least $V^\pi(\omega)$, implying that $Z(\omega)$ must be at least as large. Consider the first call ($k = 1$), and set $\bar{y}_k = A$. If $1 \in S^\pi(\omega)$, then we set $\bar{x}_{1, \bar{y}_k} = 1$, and use Constraint (3) to obtain \bar{y}_2 . If this is not the case, we set $\bar{x}_{1a} = 0$ for all a , but still use Constraint (3) to obtain \bar{y}_2 . We proceed in a similar fashion for all remaining calls, in order of increasing k ; this results in a solution that serves the same calls as policy π .

We claim that (\bar{x}, \bar{y}) satisfies $\bar{y}_k \geq A_k^\pi(\omega)$ for all calls k . This has two consequences. First, we have that (\bar{x}, \bar{y}) is feasible, as calls are never admitted when all ambulances are busy, and the constraints of (IP(ω)) are satisfied by construction. Second, we have that the bounding system collects at least as much reward from each call as the original system, as $v(\cdot)$ is nondecreasing. Combining these two implications yields $Z(\omega) \geq V^\pi(\omega)$, as desired.

It remains to verify the claim, which we do by proving a slightly stronger claim: that the service time for each call that is treated in the original system is at least as large as the corresponding service time in the bounding system. We proceed via induction on k . Consider first the base case ($k = 1$). If $1 \notin S^\pi(\omega)$, then the base case holds trivially, as the call is not served. Otherwise, $G_A^{-1}(U_k(\omega))$ lower bounds the service time of the call in the original system, by construction of the dominating service time distributions G_a .

Now suppose the claim holds until the k^{th} call arrives, where $k > 1$. The induction hypothesis implies $\bar{y}_k \geq A_k^\pi(\omega)$. The case where $k \notin S^\pi(\omega)$ is again trivial. Otherwise, $G_{A_k^\pi(\omega)}^{-1}(U_k(\omega))$ lower bounds the service time of the call in the original system. Although it may be that $\bar{y}_k > A_k^\pi(\omega)$, we have that $G_{\bar{y}_k}^{-1}(U_k(\omega)) \leq G_{A_k^\pi(\omega)}^{-1}(U_k(\omega))$, as the distributions G_a are stochastically decreasing in a by construction. This completes the induction, and we are done. \square

We have thus shown the following result:

THEOREM 1. *Let π be a redeployment policy satisfying the criteria specified in Section 2, and let Z be as defined in Section 4.1. Then $E[Z]$ is an upper bound on the expected reward collected by any policy π on the horizon $[0, T]$.*

We estimate $E[Z]$ using Monte Carlo methods, by generating and averaging i.i.d. replications of Z , and solving an instance of integer program for each replication.

EXAMPLE 2. Consider once again the system from Example 1 in Section 3.3. We observe that all of the randomness in this system is captured by the random variables U_1, \dots, U_6 . Because G_1 corresponds to a random variable taking on one of two values with probability 0.5, and G_2 is a Dirac delta function, each sample path ω can be summarized by specifying whether $U_k \leq 0.5$ or $U_k > 0.5$ for each call k . Thus, to compute $E[Z]$, we need only consider $2^6 = 64$ sample paths. Enumerating over each possibility, we find that $Z(\omega) = 3.5$ for each ω . This implies $E[Z] = 3.5$, which indeed upper bounds the performance of the heuristic policy we considered in Example 1.

5. Computational Study

In this section, we conduct numerical experiments on a variety of EMS systems (some realistic, some not realistic) to evaluate the tightness of our upper bound.

5.1. Reference Bounds

We cannot directly compare our bound to that of Maxwell et al. (2014), as we assume that the EMS operates as a loss system. However, we provide context via a lower bound obtained by simulating the performance of a heuristic redeployment policy, and an alternative upper bound stemming from an “information relaxation” of the original problem.

5.1.1. Lower Bound Our heuristic policy is adapted from that developed by Jagtenberg et al. (2015). Their policy dispatches the closest available ambulance to arriving calls, and when an ambulance becomes idle, redeploys it to the base providing the largest marginal increase in coverage, where coverage is computed as in the objective function of the MEXCLP due to Dakin (1983). We do not specify this objective function here, but note that it is a function of q , an estimate of average ambulance utilization, as well as C , the locations of idle ambulances at the time a redeployment decision is made; we write this quantity as $M(C, q)$.

We make two modifications to this policy to improve its performance. First, when computing coverage, we consider ambulances undergoing redeployment to be idle at their destination bases. Second, we take into account the travel time associated with redeployment moves, as a short relocation yielding a moderate improvement in coverage may be preferable to a long relocation achieving a slightly higher increase in coverage. Letting $\Delta M_j(C, q)$ denote the marginal increase in coverage associated with placing an additional ambulance at base j , i the location of the ambulance to be redeployed, and $\alpha \geq 0$, we let

$$V_{ij} := \frac{\Delta M_j(C, q)}{\max\{(t_{ij})^\alpha, 1\}}$$

denote the “value” associated with a redeployment move from node i to base j . Our heuristic policy relocates the ambulance to the base attaining the largest V_{ij} value. We use a maximum in the denominator to avoid potential divisions by zero, and to avoid excessively inflating the value of short redeployments. This is a generalization of the metric used by Jagtenberg et al., which we can recover by setting $\alpha = 0$. We can generalize their policy further by treating q as a parameter that can be used to tune the redeployment policy, rather than as measure of system-wide utilization. This allows us to improve our lower bound via a two-dimensional grid search over α and q .

5.1.2. Information-Based Upper Bound Our alternative upper bound stems from a perfect information relaxation of the original problem. This entails solving a problem in which the decision-maker knows the arrival times, locations, and service requirements associated with every call at the start of the horizon— quantities that are typically not known before calls arrive. The optimal objective value associated with this relaxed problem, in expectation, upper bounds the performance of any policy in the original problem; this follows by formulating the problem described in Section 2 as a stochastic dynamic program, and using results from Brown et al. (2010); we refer the reader to Chong (2016) for further details.

If we assume, contrary to our problem formulation in Section 2, that the decision-maker can reject arriving calls when ambulances are available, then we can solve the relaxed problem for a given sample path by instead solving an equivalent linear program that is similar in spirit to the integer program appearing in Yue et al. (2012); we refer the reader to Chong (2016) for the formulation, and for the corresponding proofs. This additional assumption loosens the upper bound, but is nonetheless valid, and provides us with another way to evaluate our upper bound from Section 4.

5.2. Findings

5.2.1. A First Example. We begin by considering a small hypothetical EMS in which our upper bound is fairly tight. The service area, which is illustrated in Figure 1, is a 9 mile \times 9 mile square region, which we divide into 81 square cells. We treat each cell as a demand node, and define the distance between any two nodes to be the Manhattan distance between the centroids of the corresponding cells. Calls arrive according to a homogeneous Poisson process with rate λ . To treat these calls, the EMS provider operates 4 bases and deploys 4 ambulances, each of which travel at a constant speed of 30 miles per hour. We assume that the response time threshold is 9 minutes, and that the chute time is 1 minute. Thus, an ambulance can provide an adequate response if it travels no more than 4 miles to reach a call. Service times follow a Weibull distribution with a scale parameter of 30 and a shape parameter of 3, which has a mean of 26.8 minutes and a standard deviation of 9.7 minutes.

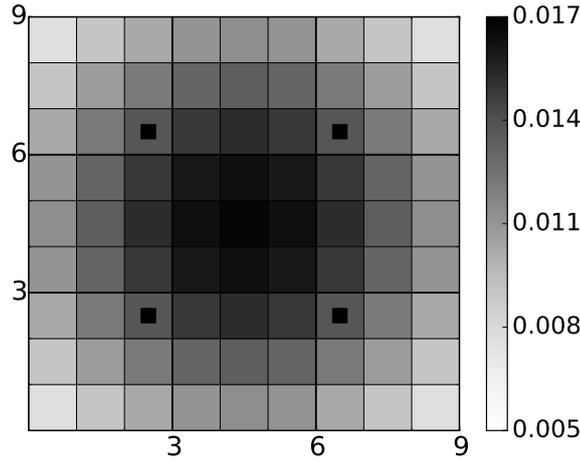


Figure 1 Distribution of demand and base locations (squares) in our hypothetical EMS system. Darker cells indicate areas with higher demand (larger d_i -values).

Table 1 Estimated lower and upper bounds on the proportion of timely responses in a 24-hour period, for two choices of arrival rates. Half-widths of 95% confidence intervals are at most 0.001.

	$\lambda = 2$ calls/hour	$\lambda = 3$ calls/hour
Lower Bound	0.842	0.741
Upper Bound	0.862	0.780
Perfect Information Upper Bound	0.976	0.935

We compute the bounding function v by solving the appropriate MCLP instances to optimality. Similarly, we compute the bounding service time distributions $\{G_a : 1 \leq a \leq A\}$ by solving the corresponding p -median problems for times ranging from 1 to 90 minutes in increments of 1 minute, also to within 0.5% of optimality. We compute our lower and upper bounds by taking averaging 1000 i.i.d. replications, each with a length of 24 hours, for the cases when $\lambda = 2$ and 3 calls per hour (representing situations in which the system is lightly and moderately loaded). To obtain our lower bound, we simulate our heuristic policy on the same sample paths for values of α and q ranging from 0.05 to 0.95 in increments of 0.05, and store the objective value attained by the best policy. We summarize our results in Table 1.

In this example, the gap between our lower bound and our upper bound is narrow; our heuristic policy adequately responds to within 4% of the calls that can be reached under the optimal redeployment policy. By contrast, the perfect information upper bound is loose, which is somewhat striking, as the systems are not heavily loaded; average ambulance utilization under the heuristic

policy is roughly 0.26 when $\lambda = 2$, and 0.39 when $\lambda = 3$. This suggests that there is significant value in learning call locations in advance, as is the case in the perfect information upper bound—which is intuitive, as this information allows the decision-maker to position ambulances closer to arriving calls, thus increasing the likelihood of timely responses.

5.2.2. Realistic Models. Next, we study two EMS systems loosely modeled after those operated in Edmonton, Canada and Melbourne, Australia. Our models are realistic in that we use transportation networks based upon simplified versions of the road networks in both cities, and we obtain input parameters by drawing from the same datasets used by Maxwell et al. (2014), which reflect real data. However, we use the term “loosely” because we make simplifying assumptions—specifically, that system dynamics (including the arrival process, edge travel times, and the size of the fleet) are time-homogeneous, and that calls do not require transport to a hospital. Thus, we do not intend for the results below to reflect how the systems in Edmonton and Melbourne actually perform. Rather, we are simply using the data available to us to gain insight into how our upper bound performs on larger, more realistic systems.

Our graph of Edmonton, which is identical to that used in Maxwell et al. (2014), contains 4,725 demand nodes, which are served by 16 ambulances distributed among 11 bases. Our graph of the city of Melbourne (also from Maxwell et al. (2014)) features 1,413 demand nodes, 97 ambulances, and 87 bases. As in the Section 5.2.1, arrivals occur according to a homogeneous Poisson process, and we use the same response time threshold, chute time, and service time distribution.

We again compute the function v by solving MCLP instances to optimality, but we compute the distributions $\{G_a : 1 \leq a \leq A\}$ by solving the linear programming relaxations of the corresponding p -median problems (due to their size). In doing so, we do not invalidate the upper bound, but we loosen it slightly. However, computational experiments suggest that the integrality gap is fairly small. To compute our lower and upper bounds, we again use 1000 replications, each with a length of 24 hours; computational experiments indicate that a heuristic policy where $\alpha = 0.50$ and $q = 0.70$ works well in Edmonton, and a heuristic policy where $\alpha = 0.85$ and $q = 0.75$ works well in Melbourne. We summarize our results in Table 2.

Table 2 Estimated lower and upper bounds in systems loosely modeled after Edmonton and Melbourne EMS. Half-widths of 95% confidence intervals are at most 0.001.

	Edmonton		Melbourne	
	$\lambda = 6/\text{hour}$	$\lambda = 8/\text{hour}$	$\lambda = 32/\text{hour}$	$\lambda = 40/\text{hour}$
Lower Bound	0.805	0.771	0.844	0.824
Upper Bound	0.854	0.848	0.888	0.888
Perfect Information Upper Bound	0.856	0.856	0.887	0.887

Our upper bounds appear to be insensitive to the level of congestion in either system. This may seem counterintuitive, as the numbers in Table 2 suggest that even in situations when the decision-maker has perfect information, a significant proportion of arrivals are not being adequately treated. However, these results can be explained by the fact that in both cities, a nontrivial fraction of demand originates from locations that cannot be reached by any base within the response time threshold— 0.144 in Edmonton and 0.112 in Melbourne. Thus, our upper bounds are close to the theoretical maximum possible performance that can be attained in either system.

Although the gap between our lower and upper bounds is not large, our upper bound does not improve upon the theoretical upper bound, or upon the perfect information upper bound. To better understand the discrepancy we observe between the results in this section and those in Section 5.2.1, we conduct a third set of experiments on a range of EMS systems with varying characteristics.

5.2.3. Exploratory Models. Figure 2 illustrates the 9 hypothetical EMS systems we consider. They are adapted from the “artificial cities” considered in Maxwell et al. (2014). The systems vary by the number of demand modes (1, 2, or 5), as well as by the dispersion of demand within the service area (low, medium, or high). The EMS provider in each system responds to calls within a 15 mile \times 15 mile square region, which is divided into 225 square cells, and operates 25 bases. As in Section 5.2.1, we compute distances using the Manhattan metric, and the same service time distribution, chute time, response time threshold, and ambulance travel speed. We consider fleets operating 7 and 10 ambulances, with arrival rates of 3 and 6 per hour, respectively. Table 3 lists the bounds we obtain in our 9 systems, averaged over 1000 i.i.d. replications.

The perfect information upper bound is loose in each of our experiments, and appears to be insensitive to how demand is distributed within the system. This is intuitive, as when utilization is

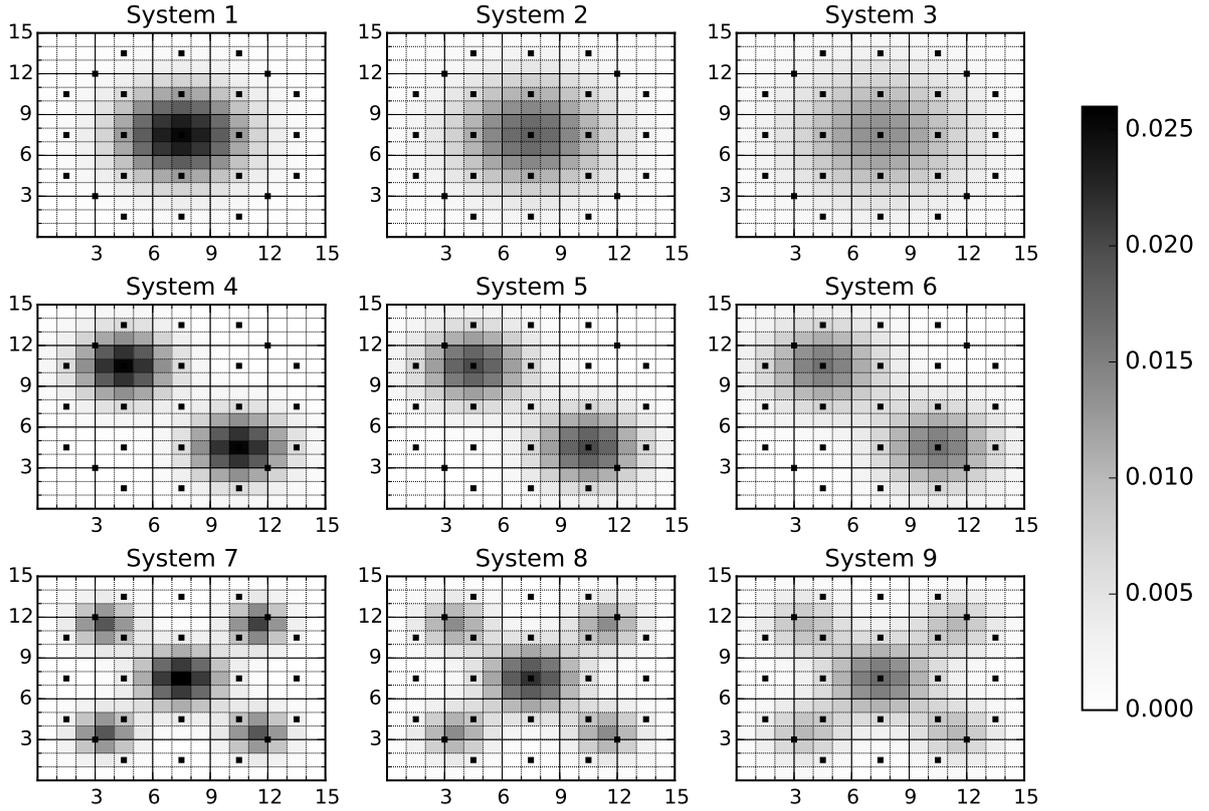


Figure 2 Distribution of demand and base locations in our 9 hypothetical EMS systems. Darker cells indicate areas with higher demand (larger d_i -values).

Table 3 Estimates of our lower and upper bounds in 9 hypothetical EMS systems. Once again, half-widths of 95% confidence intervals are at most 0.001.

7 ambulances, $\lambda = 3$	Sys 1	Sys 2	Sys 3	Sys 4	Sys 5	Sys 6	Sys 7	Sys 8	Sys 9
Lower Bound	0.918	0.869	0.825	0.928	0.852	0.799	0.859	0.853	0.836
Upper Bound	0.972	0.949	0.926	0.979	0.967	0.953	0.960	0.956	0.946
Pf Info Upper Bound	0.996	0.993	0.990	0.996	0.991	0.985	0.995	0.995	0.993
10 ambulances, $\lambda = 6$									
Lower Bound	0.939	0.895	0.857	0.919	0.898	0.880	0.878	0.874	0.861
Upper Bound	0.985	0.973	0.961	0.988	0.982	0.975	0.979	0.977	0.971
Pf Info Upper Bound	0.998	0.997	0.995	0.996	0.995	0.993	0.995	0.995	0.994

low, variation in the locations of arriving calls have a smaller effect on performance when perfect information is available. The decision-maker can, to an extent, plan for future arrivals by positioning ambulances appropriately, regardless of where they may occur.

Our upper bound is tighter than the perfect information bound, particularly when demand is more dispersed, or when the system operates a smaller fleet. This indicates that the advantage of our upper bound may lie in its ability to take into account uncertainty in demand, and in more

accurately modeling the performance of the system during periods of congestion. Nonetheless, the gap between our lower and upper bounds is significant.

We conjecture this may be due to the upper bound being loose, rather than to the heuristic policy being poor. This is because the bounding functions v and $\{G_a : 1 \leq a \leq A\}$ were developed by making the optimistic assumption that at any given time, ambulances are “optimally positioned”—that is, in a way that maximizes the likelihood of timely responses and short service times. When the system becomes congested, ambulances may be in a suboptimal configuration, resulting in the probability of a timely response being smaller than the function v would suggest. Furthermore, the distribution function G_a may underestimate the travel time of the responding ambulance, as dispatches may need to be made from distant bases when the system is congested. These effects may be more pronounced in large-scale systems.

This conjecture is consistent with our previous results: our bound performs well in the small-scale system in Section 5.2.1, but struggles with the larger-scale models in Section 5.2.2. Taken together, our experiments suggest that our upper bound is best suited for small-scale systems, but tends to improve upon the perfect information bound even when used outside of this setting.

6. Conclusion

In this paper, we constructed an upper bound on the performance of ambulance redeployment policies in loss systems. Our work builds upon that by Maxwell et al. (2014), whose bound is valid only for systems that maintain queues for calls that cannot be immediately served. The adaptation is nontrivial, and requires the introduction of new ideas. In particular, we formulate an integer program to which the set of decisions made by any nonanticipative redeployment policy (along a given sample path) correspond to a feasible solution. Computational experiments suggest that while our bound tends to be tighter than that originating from a perfect information relaxation of the ambulance redeployment problem, it is most effective in small-scale EMS systems. We conjecture that improving the gap between our lower and upper bounds hinges less upon improving the heuristic policy than upon developing refinements of the bounding functions v and $\{G_a : 1 \leq a \leq A\}$.

A natural question to ask is whether upper bounds can be obtained using other approaches. One way to do so may be to study bounds based upon relaxations of the stochastic dynamic programming formulation of the problem. This could entail, for instance, relaxing a coupling constraint, introducing Lagrange multipliers, and proceeding as in Adelman and Mersereau (2008). Alternatively, one could consider an information relaxation of the original problem. Although we consider a perfect information relaxation in our numerical study in Section 5, one might tighten the upper bound by applying an “information penalty” to decisions violating nonanticipativity constraints in the style of Brown et al. (2010). We discuss our attempts to do so in Chong (2016), but we have yet to find a penalty that meaningfully improves upon the perfect information upper bound. We also propose the construction of a more effective penalty as a direction for future research.

Acknowledgments

This work was partially supported by National Science Foundation Grant CMMI-1200315.

References

- Adelman, D., A.J. Mersereau. 2008. Relaxations of weakly coupled stochastic dynamic programs. *Operations Research* **56**(3) 712–727.
- Alanis, R., A. Ingolfsson, B. Kolfal. 2013. A Markov chain model for an EMS system with repositioning. *Production and Operations Management* **22**(1) 216–231.
- Andersson, T., P. Vaerband. 2015. Decision support tools for ambulance dispatch and relocation. *Journal of the Operational Research Society* **58** 195–201.
- Berman, O. 1981a. Dynamic repositioning of indistinguishable service units on transportation networks. *Transportation Science* **15**(2) 115–136.
- Berman, O. 1981b. Repositioning of distinguishable urban service units on networks. *Computers and Operations Research* **8**(2) 105–118.
- Berman, O. 1981c. Repositioning of two distinguishable service vehicles on networks. *IEEE Transactions on Systems, Man, and Cybernetics* **11**(3) 187–193.

- Bledsoe, B.E. 2003. Ems myth #7: System Status Management (SSM) lowers response times and enhances patient care. *Emergency Medical Services* **32**(9) 158–9.
- Brown, D.B., J. E. Smith, P. Sun. 2010. Information relaxations and duality in stochastic dynamic programs. *Operations Research* **58**(4) 785–801.
- Chong, K. 2016. Models for decision-making and performance evaluation in emergency medical service systems. Ph.D. thesis, Cornell University. In progress.
- Church, R., C. ReVelle. 1974. The maximal covering location problem. *Papers of the Regional Science Association* **32**(1) 101–118.
- Daksin, M. S. 1983. A maximum expected covering location model: Formulation, properties, and heuristic solution. *Transportation Science* **17**(1) 48–70.
- Gendreau, M., G. Laporte, F. Semet. 2001. A dynamic model and parallel tabu search heuristic for real-time ambulance relocation. *Parallel Computing* **27**(12) 1641–53.
- Gendreau, M., G. Laporte, F. Semet. 2006. The maximal expected coverage relocation problem for emergency vehicles. *Journal of the Operational Research Society* **57** 22–8.
- Henderson, S.G. 2011. Operations research tools for addressing current challenges in emergency medical services. *Wiley Encyclopedia of Operations Research and Management Science* .
- Jagtenberg, C.J., S. Bhulai, R.D. van der Mei. 2015. An efficient heuristic for real-time ambulance redeployment. *Operations Research for Health Care* **4** 27–35.
- Maxwell, M.S., E.C. Ni, C. Tong, S.G. Henderson, H. Topaloglu, S.R. Hunter. 2014. A bound on the performance of an optimal ambulance redeployment policy. *Operations Research* **62**(5) 1014–1027.
- Maxwell, M.S., M. Restrepo, S.G. Henderson, H. Topaloglu. 2010. Approximate dynamic programming for ambulance redeployment. *INFORMS Journal on Computing* **22**(2) 266–281.
- McLay, L.A., M.E. Mayorga. 2013. An optimal dispatching model for server-to-customer systems with classification errors. *IIE Transactions* **45**(1) 1–24.
- myFoxDetroit. 2011. No ambulance available for shooting victims. URL <http://www.myfoxdetroit.com/story/18473311>. Accessed July 6, 2015.

- Nair, R., E. Miller-Hooks. 2009. Evaluation of relocation strategies for emergency medical service vehicles. *Transportation Research Record: Journal of the Transportation Research Board* **2137** 63–73.
- Naoum-Sawaya, J., S. Elhedhli. 2013. A stochastic optimization model for real-time ambulance redeployment. *Computers and Operations Research* **40** 1972–1978.
- Pedersen, L. 2012. Too few paramedics to answer call: Union official. Toronto Sun, May 2012. URL <http://www.torontosun.com/2012/05/13/too-few-paramedics-to-answer-call-union-official>. Accessed July 6, 2015.
- Schmid, V. 2012. Solving the dynamic ambulance relocation and dispatching problem using approximate dynamic programming. *European Journal of Operations Research* **219**(3) 611–621.
- Sinnema, J. 2010. ER waits keep paramedics away in calls. URL <http://www.canada.com/story.html?id=c265934a-267e-40a4-a5e3-2b4872f9fd65>. Accessed July 6, 2015.
- Sinnema, J. 2012. 'Busy' ambulance system causes concern of paramedics. Edmonton Journal. URL <http://www2.canada.com/edmontonjournal/news/archives/story.html?id=0a795570-199f-4ed0-91a5-c1c562743faa>. Accessed July 6, 2015.
- Stout, J. 1989. System status management: The fact is, it's everywhere. *JEMS: A Journal of Emergency Medical Services* **14**(4) 65–71.
- Sudtachat, K., M.E. Mayorga, L.A. McLay. 2016. A nested-compliance table policy for emergency medical service systems under relocation. *Omega* **58** 154–168.
- van Barneveld, T.C. 2016. The minimum expected penalty relocation problem for the computation of compliance tables for ambulance vehicles. *INFORMS Journal on Computing*. To appear.
- van Barneveld, T.C., S. Bhulai, R.D. van der Mei. 2015. A dynamic ambulance management model for rural areas. *Healthcare Management Science*. To appear.
- van Barneveld, T.C., S. Bhulai, R.D. van der Mei. 2016. The effect of ambulance relocations on the performance of ambulance service providers. *European Journal of Operational Research*. In revision.
- Williams, D.M. 2009. JEMS 2008 200 city survey: The future is your choice. *JEMS: A Journal of Emergency Medical Services* **34**(2) 36–51.

Yue, Y., L. Marla, R. Krishnan. 2012. An efficient simulation-based approach to ambulance fleet allocation and dynamic redeployment. *AAAI Conference on Artificial Intelligence (AAAI)*.

Zhang, L. 2012. Simulation optimisation and Markov models for dynamic ambulance redeployment. Ph.D. thesis, The University of Auckland, New Zealand.