

Optimizing Call Center Staffing using Simulation and Analytic Center Cutting Plane Methods

Július Atlason, jatlason@umich.edu
Marina A. Epelman, mepelman@umich.edu

Department of Industrial and Operations Engineering
University of Michigan
Ann Arbor, MI 48109-2117, USA

Shane G. Henderson, sgh9@cornell.edu

School of Operations Research and Industrial Engineering
Cornell University
Ithaca, NY 14853, USA

August 4, 2004

Abstract

We present a simulation-based analytic center cutting plane method to solve a sample average approximation of a call center problem of minimizing staffing costs, while maintaining an acceptable level of service in multiple time periods. We establish convergence of the method when the service level functions are discrete pseudoconcave. An extensive numerical study of a moderately large call center shows that the method is robust, and, in most of the test cases, outperforms traditional staffing heuristics that are based on analytical queuing methods.

1 Introduction

The call center (also referred to as a contact center) is an important component of the operations of many organizations. A significant fraction of the cost of operating a call center is staffing cost (Gans et al., 2003). In this paper we describe a method for selecting staff (agent) levels that minimize cost while simultaneously ensuring satisfactory customer service. Simulation is used to report service level performance for a given set of staffing levels, and an analytic center cutting plane method guides the selection of staffing levels.

This problem has received a great deal of attention in the literature, and so one can reasonably ask why there is a need for a computational tool of this sort. To answer that question we first need to describe the overall staffing process. There are variations on the following theme (e.g., Castillo et al., 2003), but the essential structure is sequential in nature and is as follows (Mason et al., 1998).

1. (Forecasting) Obtain forecasts of customer load over the planning horizon, which is typically one or two weeks long. The horizon is usually broken into short periods that are typically between 15 minutes and 1 hour long.
2. (Work requirements) Determine the minimum number of agents needed during each period to ensure satisfactory customer service. Service is typically measured in terms of customer waiting times and/or abandonment rates in the queue.
3. (Shift construction) Select staff shifts that cover the requirements. This problem is usually solved through a set-covering integer program; see Mason et al. (1998) for more details.
4. (Rostering) Allocate employees to the shifts.

The focus in this paper is on Steps 2 and 3. Steps 1 and 4 are not considered further.

Step 2 is usually accomplished through the use of analytical results for simple queuing models. Green et al. (2001) coined the term SIPP (stationary, independent, period by period) to describe the general approach. In the SIPP approach, each period of the day is considered independently of other periods, the arrival process is considered to be stationary in that period, and one approximates performance in the period by a steady-state performance measure that is usually easily obtained from analytical results for particular queuing models. Heuristics are used to select input parameters for each period that yield a close match between the predictions and reality. For a wide variety of queuing models, this procedure results in some form of the “square-root rule,” which is a rule of thumb that provides surprisingly successful staffing level suggestions. See, e.g., Borst et al. (2004); Kolesar and Green (1998); Jennings et al. (1996) for more on the square-root rule.

The SIPP approach is appealing from the standpoint of computational tractability and due to the insights it provides. However, there are cases where the SIPP approach does not do as well as one might hope (Green et al., 2001, 2003). This can occur, for example, when the use of steady-state measures to represent performance over a short period is inappropriate. See Whitt (1991) for more on this point. Moreover, call centers can be somewhat complex in structure, and this complexity can make it difficult to identify a queuing model of the center that is both mathematically tractable and a reasonable match to the true system. In such cases, simulation

is a viable alternative. Indeed, simulation is now increasingly used in Step 2; see Section VIII of Mandelbaum (2003) for many examples.

Further motivation for the use of simulation involves the linkage between staffing decisions in adjacent periods. Boosting staffing levels in one period can often help in reducing workload in subsequent periods, so that there can be linkage in performance between different periods. Such linkage can imply that there are multiple solutions to the work requirements problem that can offer valuable flexibility in Step 3. Traditional queuing approaches are not satisfactory in the presence of such linkage between periods, and in such cases one turns to simulation or other numerical methods. Indeed, Green et al. (2001, 2003) solve a system of ordinary differential equations through numerical integration to get the “exact” results for their models in order to compare the performance of various heuristics.

Assuming that one uses simulation or some other numerical method to predict performance in the call center, one then needs to devise a method to guide the selection of potential staffing levels to be evaluated through simulation. There have been several suggestions in the literature, all of which explicitly capture the linkage between periods in an attempt to realize cost savings. Like Green et al. (2001, 2003), Ingolfsson et al. (2002) use numerical integration to compute service level performance for a proposed set of staffing levels, and a genetic algorithm to guide the search. Ingolfsson et al. (2003) again use a numerical method to compute service level performance, and integer programming to guide the search. Castillo et al. (2003) devise a method for randomly generating sets of staff shifts, then use simulation to evaluate the service level performance of each set of generated staff shifts, and finally, plot the cost versus service level of the potential solutions to identify an efficient frontier. Atlason et al. (2004) use simulation to evaluate service level performance of a proposed set of shifts, and use integer programming in conjunction with Kelley’s cutting plane method (Kelley, Jr., 1960) to guide the search.

The approach in Atlason et al. (2004) relies on an assumption that service in a period is concave and componentwise-increasing as a function of the staffing level vector. To understand this assumption, consider a single period problem. Increasing the staffing level should lead to improved performance. Furthermore, one might expect “diminishing returns” as the staffing level increases, so that performance would be concave in staffing level. Empirical results suggest that this intuition is correct, at least for sufficiently high staffing levels. But for low staffing levels, the empirical results suggest that performance is increasing and *convex* in the staffing level. So performance appears to follow an “S-shaped” curve (Ingolfsson et al., 2003; Atlason et al., 2004) in one dimension.

This non-concavity can cause the cutting plane method of Atlason et al. (2004) to cut off feasible solutions, and the problem can be so severe as to lead to the algorithm suggesting impractical staffing plans. Nevertheless, the ability of the Atlason et al. (2004) approach to efficiently sift through the combinatorially huge number of potential staffing plans is appealing. One might ask whether there is a similar optimization approach that can efficiently search through alternative staffing plans while satisfactorily dealing with S-shaped curves and their multidimensional extensions. That is the subject of this paper.

We combine a relaxation on the assumption of concavity, i.e., pseudoconcavity, and additional techniques to handle both the S-shaped curves alluded to above, as well as multidimensional behavior seen in examples like that depicted in Figure 2; see Section 4.1 for more details.

Assuming that the service level functions are pseudoconcave, one can then develop an analytic center cutting plane algorithm that efficiently searches the combinatorially large space of potential staffing plans. In essence the algorithm works as follows. One begins with a polyhedron that contains an optimal solution to the staffing problem. At each stage, the algorithm selects a staffing plan that lies close to the analytic center of the polyhedron and runs a simulation at that point to determine service level performance. Depending on the results of the simulation an “optimality cut” or one or more feasibility cuts are added, thereby shrinking the polyhedron. The algorithm terminates when the polyhedron contains no integer points, or when the difference between upper and lower bounds on the optimal objective is sufficiently small.

We view the contributions of this paper as follows.

1. Under realistic assumptions on the service level functions, we give an algorithm for solving to optimality the combined Step 2 - Step 3 problem. The combined procedure explicitly addresses linkage between periods.
2. We give conditions under which the algorithm provably converges.
3. We compare the proposed algorithm to what can be reasonably viewed as the current best practice to better understand the properties of the algorithm. This comparison has to take place with a model of limited complexity so that the existing methods can be applied, but is nevertheless quite illuminating.

The numerical results in Section 5 show that the analytic center cutting plane method outlined here outperforms, or at least equals, the SIPP heuristics in every case in which shift structure is explicitly considered, which is the setting we are interested in. In that sense, it is a robust procedure that can be applied in a near black-box fashion. Of course, these extremely appealing properties have to be traded off against the computational cost of the procedure, which

is not inconsiderable. We are actively considering methods for reducing the computational effort of the procedure; see Section 6.

The remainder of this paper is organized as follows. We formulate the call center staffing problem in Section 2. In Section 3 we review cutting plane methods for continuous problems with pseudoconcave constraints. Section 4 describes modifications for discrete problems, such as the call center staffing problem, and proves that the algorithm converges. We compare the proposed algorithm with the SIPP methods described in Green et al. (2001, 2003) in Section 5. Section 6 concludes the paper and discusses future research questions.

2 Problem formulation

In this section we formulate the call center staffing problem of minimizing staffing cost over a planning horizon while maintaining an acceptable level of service in each of a number of time periods, with particular emphasis on the use of simulation to estimate the service levels. We say that service is acceptable in a period if a threshold long-run percentage of calls that are required to be answered within a certain time limit τ is met. This is the usual “ π percent of calls are answered within τ seconds” requirement that is something of an industry standard.

Define the staffing level vector as $y = (y_1, \dots, y_p)^T$ where p is the number of time periods in the planning horizon, y_i is the number of agents working in period $i, i \in \{1, \dots, p\}$, and $(\cdot)^T$ denotes the transpose operation. Let N_i be the random number of calls that are received in period i , and let $S_i(y)$ be the random number of those calls that are satisfactorily handled, $i \in \{1, \dots, p\}$. We can express the service level constraint as

$$g_i(y) = ES_i(y) - \pi EN_i \geq 0, i = 1, \dots, p.$$

Here $g_i(y)$ gives the expected number of successful services over and above the required expected number of successful services in period i .

Define the cost function

$$\begin{aligned} f(y) = \min \quad & c^T x \\ \text{s.t.} \quad & Ax \geq y \\ & x \geq 0 \text{ and integer,} \end{aligned} \tag{1}$$

where x is a vector with j th component x_j giving the number of employees working the j th tour and A is the tour matrix such that $A_{ij} = 1$ if tour j includes period i and 0 otherwise. The j th

component of the cost vector c is the cost of the j th tour. The term “tour” refers to a collection of periods that a single employee works over the planning horizon. A tour must agree with all provisions of employee contracts, such as rules on meal breaks and a limit on the number of hours an agent can work over the planning horizon. All feasible tours are enumerated prior to problem formulation and stored in the tour matrix A . We assume that there are m feasible tours, so $x \in \mathbb{Z}_+^m$. We assume that every period is covered by at least one tour, i.e., for every i there is at least one j such that $A_{ij} = 1$. It then follows that (1) is feasible for any y . The value $f(y)$ gives the minimum cost set of shifts that can cover the desired work requirements vector y .

The call center staffing problem can now be formulated as

$$\begin{aligned} \min \quad & f(y) \\ \text{s.t.} \quad & g_i(y) \geq 0 \text{ for } i \in \{1, \dots, p\} \\ & y \geq 0 \text{ and integer.} \end{aligned} \tag{2}$$

Recall that we cannot compute the (vector-valued) service level function $g(y)$ exactly, and instead use simulation. The call center staffing problem is then a simulation optimization problem. As in Atlason et al. (2004) we adopt “sample-average approximation” (see Shapiro, 2003; Kleywegt et al., 2001) for solving simulation-optimization problems. This approach, specialized to our setting, is as follows. One first generates the simulation input data for n independent replications of the operations of the call center over the planning horizon. This data includes call arrival times, service times and so forth. The simulation data is fixed, and one then solves a deterministic optimization problem that chooses staffing levels so as to minimize staffing cost, while ensuring that average service *computed only over the generated realizations* is satisfactory. This problem can be solved using any convenient optimization algorithm.

Suppose that the service level functions $g_i(y)$ are estimated by the sample average $\bar{g}_i(y; n)$, where n is the sample size used. The sample average approximation of the call center staffing problem is then

$$\begin{aligned} \min \quad & f(y) \\ \text{s.t.} \quad & \bar{g}_i(y; n) \geq 0 \text{ for } i \in \{1, \dots, p\} \\ & y \geq 0 \text{ and integer.} \end{aligned} \tag{3}$$

Proposition 1 below gives conditions under which solutions to (3) converge to those of the true problem (2) as $n \rightarrow \infty$. We first make the following natural assumption about the cost vector.

Assumption 1. The cost vector c is positive and integer valued.

Assumption 1 implies that $f(y)$ is integer valued and, moreover, since c is positive the z -level set of f ,

$$\{y \geq 0 \text{ and integer} : \exists x \geq 0 \text{ and integer}, Ax \geq y, c^T x \leq z\},$$

is finite for any $z \in \mathbb{R}$.

Let Y^* denote the set of optimal solutions to the true problem (2).

Proposition 1. *Suppose that Assumption 1 holds. Suppose further that for each y , $\bar{g}_i(y; n) \rightarrow g_i(y)$ as $n \rightarrow \infty$ with probability 1. Finally, suppose that there is an optimal solution $y^* \in Y^*$ that is “strictly feasible,” i.e., $g_i(y^*) > 0$ for all $i = 1, \dots, p$. If y_n^* is an optimal solution to (3) for each n , then $y_n^* \in Y^*$ for n sufficiently large, with probability 1.*

This result is easily proved using techniques very similar to those used in Atlason et al. (2004) and so we omit the proof.

Proposition 1 establishes the validity of the sample average approximation approach in our setting. Much more can be said about the convergence of solutions of the sample average approximation problem and their objective values, but that is not the emphasis of this paper. We refer the interested reader to Atlason et al. (2004); Kleywegt et al. (2001).

3 The analytic center cutting plane method

In this section we state a version of the traditional analytic center cutting plane method (AC-CPM) to fix ideas and provide a departure point for a reader unfamiliar with cutting plane methods in continuous optimization.

There are many cutting plane methods for solving convex optimization problems, including what may be termed boundary methods, such as Kelley’s algorithm (Kelley, Jr., 1960) and its extensions (e.g., Westerlund and Pörn, 2002), and interior point methods, recently reviewed by Mitchell (2003). We will focus our attention on one of the latter, namely the ACCPM, which was first implemented in duMerle (1995), as pointed out in Elhedhli and Goffin (2003). The ACCPM has proven effective in terms of both theoretical complexity (Atkinson and Vaidya, 1995; Nesterov, 1995; Goffin et al., 1996; Mitchell, 2003) and practical performance on a variety of problems (Bahn et al., 1995; Mitchell, 2000, and other references in Mitchell, 2003). Software packages implementing the method are available (e.g., Peton and Vial, 2001).

Many versions of the ACCPM for convex feasibility and optimization problems have been

explored in the literature. The description we chose below borrows from Nesterov (1995), duMerle et al. (1998), and Mitchell (2003), as just some of the possible references.

Consider an optimization problem in the following general form:

$$\begin{aligned} \min \quad & b^T y \\ \text{s.t.} \quad & y \in Y, \end{aligned} \tag{4}$$

where $Y \subset \mathbb{R}^n$ is a convex set, and $b \in \mathbb{R}^n$. (A problem of minimizing a general convex function over a convex set can be easily represented in this form.) To simplify the presentation, assume that the set Y is bounded and has a non-empty interior.

To apply the ACCPM (or any other cutting plane algorithm), the feasible region Y needs to be described by a separation oracle. Such an oracle will, given an input $\hat{y} \in \mathbb{R}^n$, either correctly assert that $\hat{y} \in Y$, or otherwise return a non-zero vector $q \in \mathbb{R}^n$ such that

$$q^T(y - \hat{y}) \geq 0 \quad \forall y \in Y,$$

i.e., produce a *feasibility cut*. (Depending on how the set Y is described, the oracle might produce a deep or a shallow cut, which have the same form as the constraint above, but a nonzero right hand side.)

We now describe a typical iteration of the ACCPM. At the beginning of an iteration, we have available a finite upper bound z on the optimal value of (4), and a polyhedron $P = \{y \in \mathbb{R}^n : Dy \geq d\}$ that is known to contain the set Y . Here $D \in \mathbb{R}^{r \times n}$ and $d \in \mathbb{R}^r$ for some finite r . We first compute the (weighted) analytic center of the set $P \cap \{y \in \mathbb{R}^n : b^T y < z\}$ (for ease of presentation we assume that the set $P \cap \{y \in \mathbb{R}^n : b^T y < z\}$ is bounded), defined as the solution of the convex optimization problem

$$\min_y \{ -w \log(z - b^T y) - \sum_{j=1}^r \log(D_{j \cdot} y - d_j) \}, \tag{5}$$

where $D_{j \cdot}$ is the j th row of D and $w > 0$ is a weight constant that affects the convergence rate of the algorithm (see, for example, duMerle et al., 1998). The set $P \cap \{y \in \mathbb{R}^n : b^T y < z\}$ is often referred to as the *localization set*, since it contains all feasible solutions with objective function values lower than z .

Finding a solution to (5) with high degree of precision is a relatively simple task from a practical standpoint and can be done, e.g., via Newton's method. Let \hat{y} be the analytic center found. Next, the oracle is called with \hat{y} as the input. If $\hat{y} \in Y$, then, by construction, $b^T \hat{y} < z$,

The analytic center cutting plane method (ACCPM)

Initialization Start with a polyhedron $P^0 := \{y \in \mathbb{R}^n : D^0 y \geq d^0\}$ such that $Y \subset P^0$, and an upper bound z^0 . Let $w^0 > 0$, and set the iteration counter $k := 0$.

Step 1 If termination criterion is satisfied, then stop, and return y^* as a solution. Otherwise, solve problem (5) with $w = w^k$, $z = z^k$, and $P = P^k := \{y \in \mathbb{R}^n : D^k \geq d^k\}$; let y^k be the solution.

Step 2a If $y^k \in Y$ then let $z^{k+1} := b^T y^k$, $y^* := y^k$, $D^{k+1} := D^k$ and $d^{k+1} := d^k$.

Step 2b If $y^k \notin Y$ then generate one or more feasibility cuts at y^k . Update D^k and d^k to include the new constraints, and let D^{k+1} and d^{k+1} represent the new constraint set. Let $z^{k+1} := z^k$.

Step 3 Compute w^{k+1} and let $k := k + 1$. Go to Step 1.

Figure 1: The analytic center cutting plane method (ACCPM)

and the upper bound is lowered by taking $z := b^T \hat{y}$. Otherwise, if $\hat{y} \notin Y$, the oracle will produce a vector q providing a feasibility cut, which is then added to the description of the polyhedron P . The procedure is then repeated. A slightly more detailed description of the algorithm is presented in Figure 1.

Intuitively, the algorithm's efficiency stems from the fact that at each iteration the cut being added passes through the analytic center of the localization set, which is often located near a geometric center. Thus, the volume of the localization set reduces rapidly with each iteration.

Suppose the set Y is specified by $Y = \{y \in \mathbb{R}^n : g_i(y) \geq 0, i \in \{1, \dots, p\}\}$, where the functions $g_i(y)$, $i \in \{1, \dots, p\}$ are pseudoconcave, as defined below:

Definition 1. (Cf. Definition 3.5.10 in Bazaraa et al., 1993) Let $g : S \rightarrow \mathbb{R}$ be differentiable on S , where S is a nonempty open set in \mathbb{R}^n . The function g is said to be *pseudoconcave* if for any $\hat{y}, y \in S$ with $\nabla g(\hat{y})^T (y - \hat{y}) \leq 0$ we have $g(y) \leq g(\hat{y})$. Equivalently, if $g(y) > g(\hat{y})$, then $\nabla g(\hat{y})^T (y - \hat{y}) > 0$.

With Y in the above form, the feasibility cut at point $\hat{y} \notin Y$ which violates the i th constraint can be specified as

$$\nabla g_i(\hat{y})^T (y - \hat{y}) \geq 0, \tag{6}$$

since any solution y that satisfies $g_i(y) \geq 0$ also satisfies $g_i(y) > g_i(\hat{y})$.

4 A cutting plane method for discrete problems

In this section we describe how the ACCPM algorithm of Section 3 can be modified to solve the sample average approximation (3) of the call center staffing problem (2).

The most significant modification to the ACCPM for continuous problems is needed to take into account the fact that the feasible region of (3) is no longer a convex, or even connected, set, due to the integrality restriction on the variables. Cutting plane algorithms for nonlinear mixed integer programs have been explored in the past (see, for example, Westerlund and Pörn, 2002). However, in this and other similar papers it is assumed that the constraint functions are in fact differentiable functions of continuous variables; the integrality restrictions on the variables are, in a sense, exogenous. In such a setting the concept of a convex (continuous) nonlinear relaxation of the integer program is straightforward, and feasibility cuts are generated simply using the gradients of these continuous functions. In our setting, however, the service level functions and their sample average approximations are not defined at non-integer values of y , and devising their continuous extension, especially one that is easy to work with from the algorithmic point of view, is non-trivial at best. Therefore, we take a different approach in adapting the ACCPM to the discrete case.

As far as we know, the use of ACCPM as a solution method for nonlinear integer programs has not been reported, although the method has been successfully used to solve the linear relaxation subproblems in branching algorithms for integer programs (see, for example, Mitchell (2000) and Elhedhli and Goffin (2003), among many others).

In Section 4.1 we extend the notion of pseudoconcavity to functions of integer variables, and show how feasibility cuts can be generated assuming that the functions $\bar{g}_i(y; n)$, $i \in \{1, \dots, p\}$ are, in fact, discrete pseudoconcave. This leads to an ACCPM method for (mixed) integer programming problems satisfying the pseudoconcavity assumption; the algorithm is applicable to the types of problems considered in Westerlund and Pörn (2002), for example.

We also discuss whether the S -shaped form of the service level functions in the call center staffing problem is consistent with this assumption, and propose alternative feasibility cuts at points where it is violated.

The following Section 4.2 discusses other modifications of the original algorithm for the problem (3) and details of our implementation. Section 4.3 gives a proof of convergence.

4.1 Discrete pseudoconcave functions and feasibility cuts

We begin by defining the notions of a discrete convex set and a discrete pseudoconcave function. We denote the convex hull of the set C by $\text{conv}(C)$.

Definition 2. We say that the set $C \subseteq \mathbb{Z}^n$ is a *discrete convex set* if $C = \text{conv}(C) \cap \mathbb{Z}^n$, i.e, the set C equals the set of integer points in $\text{conv}(C)$.

Definition 3. Let $C \subseteq \mathbb{Z}^n$ be a discrete convex set and $g : C \rightarrow \mathbb{R}$. Then g is *discrete pseudoconcave* if for any $\hat{y} \in C$ there exists a vector $q(\hat{y}) \in \mathbb{R}^n$ such that for any $y \in C$,

$$q(\hat{y})^T(y - \hat{y}) \leq 0 \Rightarrow g(y) \leq g(\hat{y}).$$

Equivalently, if $g(y) > g(\hat{y})$, then $q(\hat{y})^T(y - \hat{y}) > 0$. We call the vector $q(\hat{y})$ a *pseudogradient* of g at \hat{y} .

In the continuous case, pseudoconcavity is a weaker property than concavity of a function; discrete pseudoconcavity can be viewed as a relaxation of the concave extensible function property defined in Murota (2003, p. 93).

If the functions $\bar{g}_i(y; n)$ in (3) are discrete pseudoconcave, then a feasibility cut at an integer point \hat{y} that violates the i th constraint can be obtained in the form

$$\bar{q}^i(\hat{y}; n)^T(y - \hat{y}) \geq \epsilon,$$

where $\bar{q}^i(\hat{y}; n)$ is the pseudogradient of $\bar{g}_i(\cdot; n)$ at \hat{y} , and $\epsilon > 0$ is sufficiently small.

Are the service level functions indeed discrete pseudoconcave? To provide an illustrative example, we computed the sample average of a service level function in period 2 of a simple two period model of a call center. Figure 2 (a) shows the number of calls answered on time in period 2 as a function of the staffing levels in period 1 and period 2. (Notice that this is equivalent to \bar{g} with $\pi = 0$). The number of servers ranges from 1 to 30 in period 1 and from 1 to 40 in period 2. We also include the contours of the function, which should at a minimum form convex sets; see Figure 2 (b). The function appears to follow a multi-dimensional extension of an S -shaped curve discussed in Section 1 (see also Ingolfsson et al., 2003; Atlason et al., 2004).

Relying on intuition derived from analyzing such S -shaped functions of continuous variables,¹ one can observe that the pseudoconcavity property is violated at very low staffing levels, although it appears to hold at staffing levels that have more than 10 servers in period 1 and more than 20 servers in period 2. The violation at low staffing levels is due to the fact that the performance function is essentially flat in this region. I.e., there are so few servers and calls are so backed up that no calls are answered on time, and adding an extra server would do little to alleviate the situation; cf. low values of y_1 and y_2 in Figure 2. It is certainly possible that we would encounter such a staffing level in the cutting plane algorithm for the sample average approximation of the

¹Differentiable functions of this shape can be characterized as *quasiconcave*, see Definition 3.5.5 in Bazaraa et al. (1993).

service level function; an attempt to compute or estimate a pseudogradient at this point would produce a zero vector. Therefore, as a feasibility cut at such a point \hat{y} , we might impose a lower bound on the number of servers in the period i , i.e., add the constraint $y_i \geq \hat{y}_i + 1$. This is not necessarily a valid feasibility cut, since the reason for calls backing up might be under-staffing in previous periods. In our implementation, if such a cut is added during the algorithm, we verify at termination that the corresponding constraint is not tight at the optimal solution found. If it is tight, then one should lower this bound and do more iterations of the cutting plane method.

Although the sample averages of the service level functions are not discrete pseudoconcave, they appear to be at least approximately discrete pseudoconcave for practical staffing levels. Furthermore, we have a strategy for dealing with areas where the function is not pseudoconcave. It seems reasonable, however, to assume that the *expected* service level function is pseudoconcave, since in expected value it would be likely that the probability of answering calls would always increase, although possibly by a very small value, when more servers are added. This would also hold for the sample average of the service level function, for a sufficiently large sample size. Therefore, when we prove the convergence results in Section 4.3 we assume that the sample averages of the service level functions are discrete pseudoconcave.

4.2 A simulation-based cutting plane method for the call center staffing problem with pseudoconcave service level functions

In this subsection we describe our modification of ACCPM for the problem (3).

At the beginning of a typical iteration, we have available a feasible solution y^* of (3), and an upper bound $z = f(y^*)$ on the optimal value of the problem. The point y^* is the best feasible solution found by the algorithm so far. We also have available a polyhedron $P = \{y \in \mathbb{R}^p : y \geq 0, Dy \geq d\}$ that is known to contain the feasible region.

Suppose y^* , z and P are as above. If y^* is not an optimal solution, then, since $f(y)$ takes on integer values, the localization set

$$\{y \in P : f(y) \leq z - 1, y \text{ integer}\},$$

is nonempty and contains all feasible solutions with objective values better than z . The localization set is empty precisely if y^* is an optimal solution. This observation provides grounds for the termination criteria we specify in our algorithm.

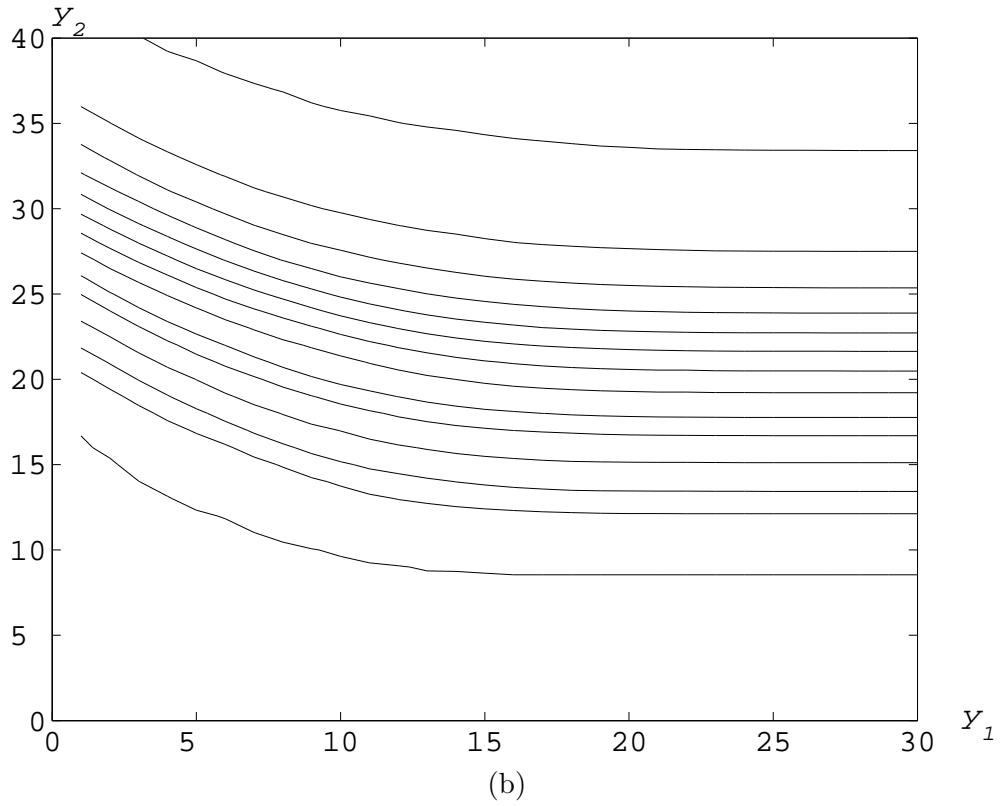
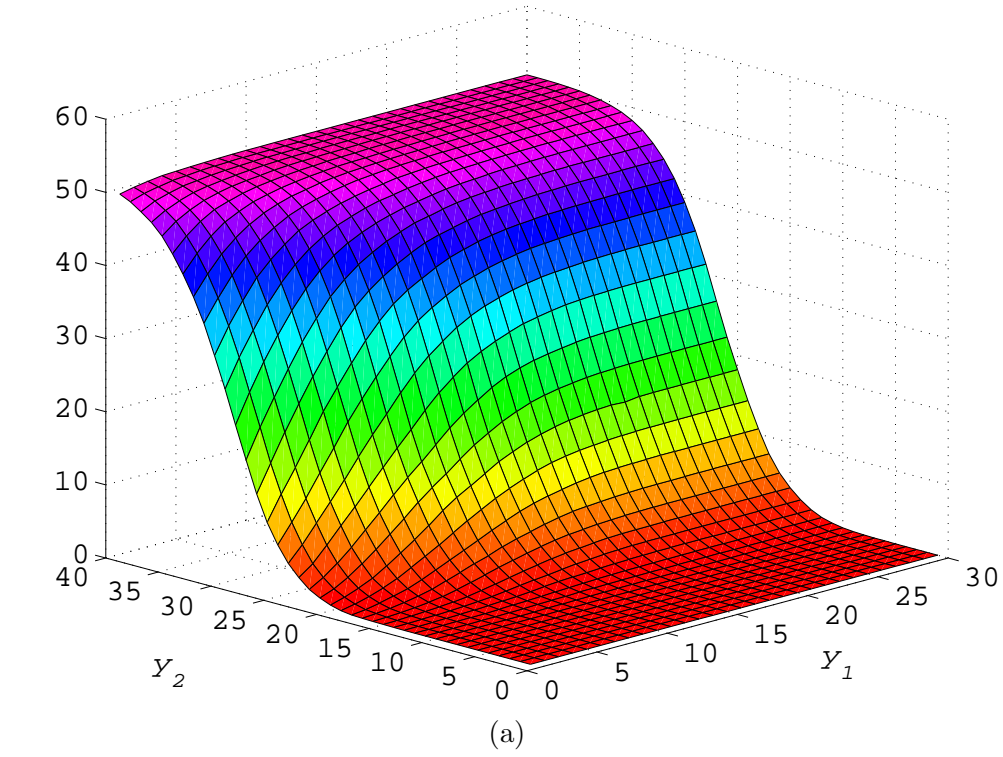


Figure 2: (a) The sample average (sample size $n = 500$) of the number of calls answered on time in period 2 as a function of the staffing levels in periods 1 and 2. (b) The contours of the service level function in (a).

Computing the next iterate. First we find the analytic center of a polyhedral relaxation of the localization set. In particular, we solve the following optimization problem:

$$\begin{aligned}
\min \quad & -w \log(z - a - c^T x) - \sum_{i=1}^p \log y_i - \sum_{j=1}^r \log(D_j \cdot y - d_j) \\
\text{s.t.} \quad & Ax \geq y \\
& x \geq 0,
\end{aligned} \tag{7}$$

where the constant $a \in (0, 1)$ ensures that only points with an objective value better than z are considered. Due to Assumption 1 the effective feasible region of (7) is bounded; hence the problem has an optimal solution so long as it has a feasible point in the effective domain of the objective. If (7) has no solution, the algorithm terminates with y^* as an optimal solution; otherwise, let $(y^{\text{ac}}, x^{\text{ac}})$ be a solution of (7). Here, $w > 0$ is the weight constant; we discuss later how this constant is determined in the algorithm.

Note that the analytic center found in the previous step is almost certainly not integer, and rounding y^{ac} to the nearest integer can result in a point outside of the localization set. To capitalize on the centrality of the analytic center in the localization set, we instead find the closest integer point in the effective feasible region of (7), i.e., solve

$$\begin{aligned}
\min \quad & \|y - y^{\text{ac}}\| \\
\text{s.t.} \quad & Ax \geq y \\
& c^T x \leq z - 1 \\
& Dy \geq d \\
& y \geq 0 \text{ and integer} \\
& x \geq 0 \text{ and integer.}
\end{aligned} \tag{8}$$

If (8) is infeasible, the algorithm terminates with y^* as an optimal solution; otherwise, let (\hat{y}, \hat{x}) be the solution of (8) and choose \hat{y} as the next iterate. Here, $\|y - y^{\text{ac}}\|$ is a measure of the distance between y and y^{ac} . If we choose the L_1 -norm as the measure, i.e., $\|y - y^{\text{ac}}\| = \sum_{i=1}^p |y_i - y_i^{\text{ac}}|$, then (8) is a linear integer program. We discuss the computational requirements of solving this problem at each iteration when we talk about the overall computational effort in relation to the computational experiments.

Estimating the service levels. Next we compute $\bar{g}_i(\hat{y}; n)$ for all i via simulation.

Adding an optimality cut. If $\bar{g}_i(\hat{y}; n) \geq 0$ for all i , then \hat{y} satisfies the service level requirements. Since $c^T \hat{x} \leq z - 1$, \hat{y} is contained in the localization set, i.e., it is the best staffing

level so far. Note that $c^T \hat{x}$ is not necessarily the cost associated with staffing level \hat{y} , since $c^T x$ is not being minimized in (8). To compute the cost associated with \hat{y} we instead solve (1) to get $f(\hat{y})$ and update $z := f(\hat{y})$.

Adding a feasibility cut. If $\bar{g}_i(\hat{y}; n) < 0$ for some i , then we add a feasibility cut for each i such that $\bar{g}_i(\hat{y}; n) < 0$. In particular, we estimate a pseudogradient, $\bar{q}^i(\hat{y}; n)$, of $\bar{g}_i(\cdot; n)$ at \hat{y} (we will discuss techniques for estimating the pseudogadients in Section 5.2). If $\bar{q}^i(\hat{y}; n) \neq 0$, we add a feasibility cut of the form

$$\bar{q}^i(\hat{y}; n)^T y \geq \bar{q}^i(\hat{y}; n)^T \hat{y} + \epsilon \quad (9)$$

for some small constant $\epsilon > 0$ (we discuss the role of ϵ in the discussion of the convergence of the method in Section 4.3). If $\bar{q}^i(\hat{y}; n) = 0$, the feasibility cut takes the form of a lower bound on the number of servers (see discussion in Section 4.1). We update D and d to reflect the cuts added.

The above procedure is then repeated. An illustration of the localization set and the feasible regions and solutions of problems (7) and (8) in each iteration is shown in Figure 3. We summarize the simulation-based analytic center cutting plane method (SACCPM) for the call center staffing problem in Figure 4. (To shorten the presentation, the description of the algorithm in Figure 4 is only specified for the case when the constraint functions $\bar{g}_i(\cdot; n)$ are in fact discrete pseudoconcave.)

The weight parameter w can be increased to “push” the weighted analytic center away from the optimality cuts. There are no theoretical results on how to choose the weights, but some computational experiments have been done to test different values of the weights, and in fact, weights on the feasibility cuts (Goffin et al., 1992; duMerle et al., 1998). The choice of the weights is problem and data dependent (see Section 5.2 for the particular choice used in our implementation).

In practice it is useful to maintain a lower bound on the optimal value of the problem (3) throughout the algorithm, in addition to the upper bound z . In particular, the algorithm can be terminated as soon as the gap between the upper and lower bounds is sufficiently small, indicating that the current “incumbent” y^* is a sufficiently good solution of the problem. A

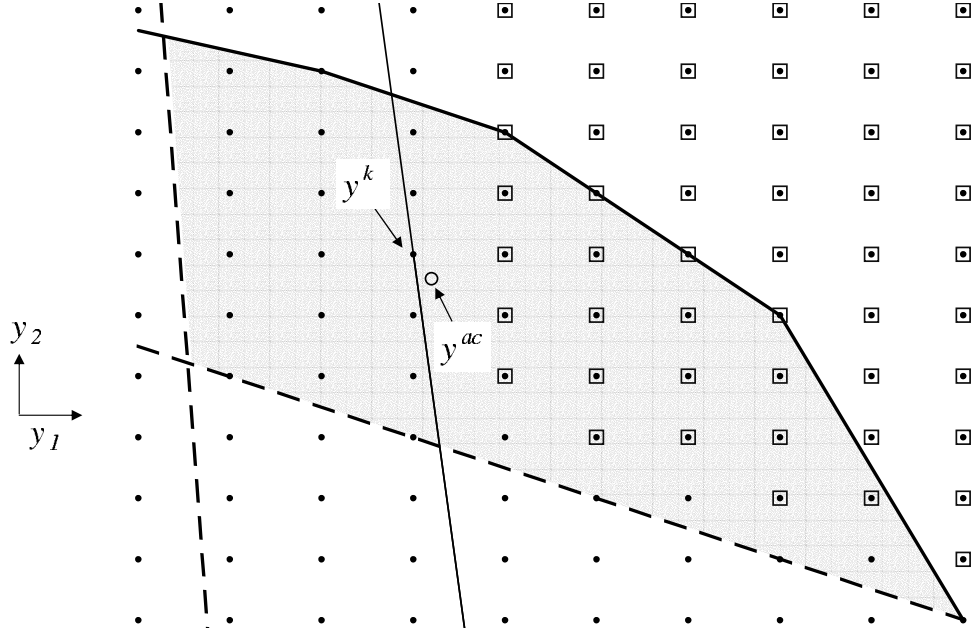


Figure 3: Illustration of the feasible region and an iterate in the SACCPM. The picture shows points in the space of the y variables in the case $p = 2$. The shaded area represents the localization set $\{y \in P : f(y) \leq z - 1, y\}$ (ignoring integer constraints). The squares represent the points that are feasible for the sample average approximation problem (3). The thick solid line segments represent an optimality cut $f(y) \leq z - 1$ and the dotted lines represent the feasibility constraints $Dy \geq d$. The point y^{ac} is the solution of the analytic center problem and y^k is the closest integer in the localization set. The thinner solid line represents a feasibility cut that would be added in this iteration.

The simulation-based analytic center cutting plane method (SACCPM)

Initialization Start with a feasible solution y^0 to (3). Let $z^0 := f(y^0)$, and $y^* := y^0$. Let $P^0 = \{y \geq 0 : D^0 y \geq d^0\}$, where D^0 and d^0 are empty. Choose an $\epsilon > 0$, $a \in (0, 1)$ and $w^0 > 0$. Let $k := 0$.

Step 1 Solve problem (7) with $w = w^k$, $z = z^k$, and $P = P^k$. If (7) does not have a (feasible) solution, then terminate with y^* as the optimal solution and z^k as the optimal value; otherwise let y^{ac} be the solution of (7).

Step 2 Solve problem (8) with $z = z^k$ and $P = P^k$. If (8) is infeasible, then terminate with y^* as the optimal solution and z^k as the optimal value; otherwise let y^k be the optimal solution of (8) found.

Step 3a If $\bar{g}_i(y^k; n) \geq 0 \forall i \in \{1, \dots, p\}$, let $z^{k+1} := f(y^k)$, $y^* := y^k$, $D^{k+1} := D^k$ and $d^{k+1} := d^k$.

Step 3b If $\bar{g}_i(y^k; n) < 0$ for some $i \in \{1, \dots, p\}$, then add the constraint $\bar{q}^i(y^k; n)^T y \geq \bar{q}^i(y^k; n)^T y^k + \epsilon$ for each i such that $\bar{g}_i(y^k; n) < 0$. Update D^k and d^k to include the added inequalities, and let $P^{k+1} = \{y \geq 0 : D^{k+1} y \geq d^{k+1}\}$ represent the new constraint set. Let $z^{k+1} := z^k$.

Step 4 Compute w^{k+1} and let $k := k + 1$. Go to Step 1.

Figure 4: The simulation-based analytic center cutting plane method (SACCPM)

lower bound can be found as the optimal objective value of

$$\begin{aligned}
 \min \quad & c^T x \\
 \text{s.t.} \quad & Ax \geq y \\
 & c^T x \leq z \\
 & Dy \geq d \\
 & y \geq 0 \text{ and integer} \\
 & x \geq 0 \text{ and integer,}
 \end{aligned} \tag{10}$$

or, for a weaker but easier to compute bound, of the linear programming relaxation of (10).

4.3 Convergence of the SACCPM

In this section we give conditions under which the SACCPM terminates with y^* as an optimal solution of (3). First, we argue that the algorithm does indeed terminate and then we show that y^* is an optimal solution of (3) at termination. To prove the results we make the following two assumptions.

Assumption 2. The functions $\bar{g}_i(y; n)$ are discrete pseudoconcave in y for all $i \in \{1, \dots, p\}$.

Assumption 3. In the implementation of the SACCPM, $\bar{q}^i(\hat{y}; n)$ is a pseudogradient of $\bar{g}_i(\cdot; n)$ at \hat{y} .

We also define the sets

$$\begin{aligned}\Gamma &:= \{y \geq 0 \text{ and integer} : f(y) \leq z^0\}, \\ \Psi &:= \Gamma \cap \{y : \bar{g}_i(y; n) \geq 0 \forall i \in \{1, \dots, p\}\}, \\ \Upsilon &:= \Gamma \setminus \Psi \text{ and} \\ I(y) &:= \{i : \bar{g}_i(y; n) < 0.\}\end{aligned}$$

In words, Γ is the set of points that are potentially visited by the algorithm and contains the set of optimal solutions. The set Ψ is the set of points in Γ that are feasible for the sample average approximation problem (3). The set Υ is the set of points in Γ that are not feasible for the sample average approximation problem (3). The set $I(y)$ is the set of periods in which the sample average of the service level function is not acceptable given the staffing levels y . The following lemma says that all solutions in Ψ satisfy potential feasibility cuts (9) for some appropriately chosen ϵ .

Lemma 2. *Let $\bar{q}^i(\hat{y}; n)$ be a pseudogradient of $\bar{g}_i(\cdot; n)$ at \hat{y} and suppose Assumptions 1 and 2 hold. Then there exists an $\tilde{\epsilon} > 0$ such that $\bar{q}^i(\hat{y}; n)^T(y - \hat{y}) \geq \tilde{\epsilon} \forall y \in \Psi, \hat{y} \in \Upsilon, i \in I(\hat{y})$.*

Proof: Let $\hat{y} \in \Upsilon$ be fixed. Suppose that $\bar{q}^i(\hat{y}; n)^T(y - \hat{y}) \leq 0$ for some $y \in \Psi$ and $i \in I(\hat{y})$. Then $\bar{g}_i(y; n) \leq \bar{g}_i(\hat{y}; n) < 0$, where the first inequality follows by Assumption 2 and the second inequality follows since $i \in I(\hat{y})$. This is a contradiction since $y \in \Psi$, and therefore $\bar{q}^i(\hat{y}; n)^T(y - \hat{y}) > 0 \forall y \in \Psi, i \in I(\hat{y})$.

Let $\epsilon(\hat{y}) = \min_{i \in I(\hat{y})} \min_{y \in \Psi} \bar{q}^i(\hat{y}; n)^T(y - \hat{y})$. The set Ψ in the inner minimum is finite by Assumption 1 and therefore the inner minimum is attained for some $y \in \Psi$. Since $I(\hat{y})$ is also finite, the outer minimum is also attained for some $i \in I(\hat{y})$. Therefore $\epsilon(\hat{y}) > 0$.

Finally, let $\tilde{\epsilon} = \min_{\hat{y} \in \Upsilon} \epsilon(\hat{y})$. Then $\tilde{\epsilon} > 0$ since Υ is finite. □

Theorem 3. *Suppose (3) has an optimal solution and that Assumptions 1, 2 and 3 hold. Furthermore, let $0 < \epsilon \leq \tilde{\epsilon}$, where $\tilde{\epsilon}$ is as in Lemma 2. Then the SACCPM terminates in a finite number of iterations returning y^* , which is an optimal solution of (3).*

Proof: The SACCPM only visits points that are feasible solutions of (8). The set of feasible solutions of (8) at every iteration is a subset of Γ and is therefore a finite set by Assumption 1. Suppose that the k th iterate y^k is the same as a previous iterate y^j for some $j < k$. If y^j

is in Ψ then $f(y^k) \leq z - 1 \leq f(y^j) - 1$, where the second inequality follows since z is the cost of the best feasible solution visited before iteration k . This is a contradiction, so y^j is not in Ψ . On the other hand, if y^j is in Υ then $\bar{g}_i(y^j; n) < 0$ for some $i \in I(y^j)$. Since y^k is in P , $\bar{q}^i(y^k; n)^T(y^k - y^j) \geq \epsilon > 0$ which is also a contradiction and y^k can therefore not be equal to y^j for any $j < k$. Therefore, the SACCPM does not visit any point in Γ more than once. Since Γ is a finite set, the algorithm is finite.

By Lemma 2 the feasibility cuts never cut off any of the feasible solutions of (3). When at termination problems (7) or (8) do not have a feasible solution it means that all the feasible solutions for (3) have an objective value greater than $z - 1$. But y^* is feasible for (3) and has an objective value of z and is, therefore, optimal, since $f(y)$ is integer by Assumption 1. \square

The theorem says that there exists an $\tilde{\epsilon}$ such that the algorithm terminates with an optimal solution of (3) if $0 < \epsilon \leq \tilde{\epsilon}$. In practice, the value of $\tilde{\epsilon}$ is unknown, but in general a “small” value for ϵ should be chosen, as in Section 5.4.

5 Numerical results

In this section we give an implementation of the SACCPM for a call center with time varying arrival rates. In Section 5.1 we describe the example that we use for the numerical experiments. In Section 5.2 we discuss how to compute, or estimate, the pseudogradients, and we also describe how the algorithm was implemented. We compare the results with staffing levels obtained by analytical queuing methods (see Section 1 for references). Analytical queuing methods based on the Erlang C formula (12) are widely used to determine staffing levels in call centers (Cleveland and Mayben, 1997). In Green et al. (2001, 2003) several heuristics for improving the performance of the basic Erlang C model are evaluated and we give a summary of these heuristics is given in Section 5.3. We believe that the methods in Green et al. (2001, 2003) are among the best analytical methods for determining required staffing levels to answer a minimum fraction of calls before a threshold time limit. Therefore we use these methods as benchmarks for our method. The results of our experiments are in Section 5.4 and we comment on the computational requirements of the SACCPM in Section 5.5.

5.1 Example

Our test model is similar to the models used in Green et al. (2001, 2003), which are call centers that can be modeled as $M(t)/M/s(t)$ queuing systems. In the Green et al. (2001) paper a call

center operating 24 hours a day and 7 days a week is studied, while the subject of the Green et al. (2003) paper are call centers with limited hours of operation. We consider a call center with limited hours of operation that is open from 6am to 12am. The call center has the following additional characteristics.

- The planning horizon consists of a single day's operation. The 18 hour planning horizon is broken into 72 time periods, each of length 15 minutes.
- In each period 80% of calls should be answered immediately. This is equivalent to setting $\tau = 0$ and $\pi = 0.8$.
- The service times for each call can be modeled as independent exponential random variables at rate μ .
- The average load over the 18 hours is R . The load is the average arrival rate divided by the average service rate and is an indication of the size of the system.
- The arrival process on any given day is a nonhomogeneous Poisson process with arrival rate at the end of period i given by $\lambda_i = \lambda(1 + \theta \sin(2\pi(t_i - 6)/18))$, where t_i is the end time of period i and is measured in the hour of the day (e.g., $t_0 = 6$ and $t_{72} = 24$). The average daily arrival rate is $\lambda = R\mu$ and θ is the relative amplitude. The arrival rate at time t such that $t_{i-1} < t \leq t_i$ is given by linear interpolation of the rates λ_{i-1} and λ_i .
- Calls are answered in the order they are received.
- When there is a reduction in the number of servers at the end of a period, any server completes a call that is already in service. A new call cannot enter service until there are fewer calls in service than there are servers for the new period.

We study the performance both in the presence of shift constraints and when there are no shift constraints.

- When there are no shift constraints, the tour matrix A is the identity matrix in \mathbb{R}^p , and the cost for each tour is equal to 1 man-period, i.e., the total staffing cost equals the total number of servers over all 72 periods.
- We assume that the shift constraints are such that each tour covers 6 hours, or 24 periods.
- The tours can only start on the hour and no later than 6 hours before the end of the day. This results in 13 tours: 6am-12pm, 7am-1pm, \dots , 5pm-11pm and 6pm-12am.
- The cost for each tour is equal to 24 man-periods per tour.

Note that we are yet to specify a value for μ , R and θ . In the experiments, we study how the SACCPM performs under two different settings (high and low) of each of these parameters.

Parameter	Low	High
μ	4 calls/hour	16 calls/hour
R	8	32
θ	.25	.75
Shifts	Yes	No

Table 1: The parameter settings in the experiments.

Along with the two settings for the shift constraints this results in a total of $2^4 = 16$ experiments. The high and low values for each parameter are given in Table 1. Green et al. (2001, 2003) additionally studied the effect of the target probability, π . The target level did not appear to have as significant an effect on the reliability of each method as the other parameters did, so we do not include different settings of it in our study. Instead of shifts, they included a factor that they call the planning period which is similar to the shift factor in that the number of servers is constant in each planning period (which can be longer than what we, and they, call a period and measure the performance in).

5.2 Implementation of the SACCPM and estimating the pseudogradients

Before we describe the implementation of the continuous and discrete optimization problems solved in each iteration of the SACCPM, we discuss what is perhaps the most challenging part of the implementation of the algorithm: estimating the pseudogradients. Up until now we have assumed that in the implementation of the SACCPM, $\bar{q}^i(\hat{y}; n)$ is a pseudogradient of $\bar{g}_i(\cdot; n)$ at \hat{y} . What makes computing, or estimating, such a pseudogradient particularly challenging is that the service level function is a discrete function of the number of agents, so a gradient, which could otherwise have been used as a pseudogradient, does not exist. In addition, we do not have a closed form expression of the function.

The simplest and perhaps the most intuitive method for estimating a gradient (or pseudogradient) when an expression for the function is unknown is the method of finite differences (see, e.g., Andradóttir, 1998). The method of finite differences can easily be extended to discrete functions, i.e., we let

$$\bar{q}_j^i(\hat{y}; n) = \bar{g}_i(\hat{y} + e_j; n) - \bar{g}_i(\hat{y}; n) \quad \forall j \in \{1, \dots, p\}, \quad (11)$$

where e_j is the j th unit vector in \mathbb{R}^p , be the estimate of the pseudogradient of $\bar{g}_i(y; n)$ at the

staffing level \bar{y} . It is clear that $p + 1$ function evaluations of \bar{g}_i are required to compute \bar{q}^i . We note, however, that estimates of the pseudogradients of the service level functions in *all* p periods can be obtained from those same $p + 1$ simulations.

In real applications of the call center staffing problem, the number of periods p can be quite large. It would therefore greatly reduce the computation time of the SACCPM if a pseudogradient can be estimated from a single simulation. The two most prominent gradient estimation techniques in the simulation literature are the likelihood ratio (LR) gradient estimation method² (see, for example, Glynn, 1990; Rubenstein and Shapiro, 1993; L'Ecuyer, 1990, 1995) and infinitesimal perturbation analysis (IPA) (see, for example, Glasserman, 1991; Fu and Hu, 1997).

The LR method is intended to estimate a gradient from a single simulation by differentiating the elements, namely the densities and the integrand, of the expected value of the service level function. The idea of IPA is to model quantities of the sample path, such as service times, as functions of the dependent variable. Under appropriate conditions the gradient of the expected value of the performance measure can be computed as the expected value of the gradient of the quantities in the simulation with respect to the dependent variable. In our formulation of the call center staffing problem, the variables of interest are the staffing levels, and are discrete. In order to to apply both the LR method and IPA, the service level functions must first be extended to a continuous domain. One such extension is to approximate the discrete service level functions by functions of service rates. Then a pseudogradient can be estimated by estimating a gradient of the service level functions with respect to the continuous rate variables.

Atlason et al. (2003) and Atlason (2004) contain a more detailed description and comparisons via numerical experiments of the three different techniques. The results indicate that, although the LR method and IPA are more tractable from a computational standpoint, the resulting gradient estimates are not as reliable as those obtained by the method of finite differences and are therefore not considered here further. In fact Atlason (2004) includes an implementation of the SACCPM using IPA pseudogradients and the the quality of the solutions from the algorithm are inferior to solutions that the algorithm gives when finite differences are used. Therefore, we used (11) to estimate the pseudogradients.

Other parts of the SACCPM were implemented as follows.

- We built a simulation model using the ProModel simulation software to compute the sample average of the service level function.
- We used Visual Basic for Applications and Microsoft Excel to store the data and to com-

²The likelihood ratio gradient estimation method is also known as the score function method.

pute the cuts.

- We used the AMPL modeling language to model and call a solver for the analytic center problem (7) and for the IP (8) of finding an integer point close to the analytic center.
- We used the MINOS solver to solve the analytic center problem (7).
- We used the CPLEX solver to solve the IP (8).
- We used the Excel solver to compute the cost $f(\hat{y})$ of a particular staffing level \hat{y} .

Finally, we describe the settings of the parameters that are specific to the SACCPM.

- We let $w^k = r$ where r is the number of feasibility cuts that have been added in iterations 1 through $k - 1$ (we let $w^k = 1$ if no feasibility cuts have been added).
- We used $\epsilon = 10^{-5}$.
- We chose $a = 1 - \epsilon$.
- Instead of a feasible starting point y^0 we started with an upper bound on the staffing levels, i.e., we added the term $-\sum_{i=1}^p \log(130 - y_i)$ to the objective of the analytic center problem (7) and the constraints $y_i \leq 130 \forall i \in \{1, \dots, p\}$ to the IP (8).
- We used a sample size of $n = 100$.

5.3 Analytical queuing methods

In this section we describe the queuing methods in Green et al. (2001, 2003) that we use as benchmarks for our computational study. Traditional queuing methods for staffing call centers assume that the process is in steady state, i.e., the arrival rate is fixed and the call center has been open for long enough that the initial state of the call center does not matter. Then assuming that the call center can be modeled as an $M/M/s$ queuing system, the probability that a customer has to wait more than τ time units, as a function of the number of servers s , is given by (Gross and Harris, 1998, p. 72)

$$P(s) = \left(\frac{\sum_{n=0}^{s-1} \frac{R^n}{n!}}{\sum_{n=0}^{s-1} \frac{R^n}{n!} + \frac{R^s}{s!(1-\rho)}} \right) e^{-(s\mu-\lambda)\tau} \quad (12)$$

where $R = \lambda/\mu$ is the load and $\rho = \lambda/s\mu$ is the server utilization. This equation is only defined for $\rho < 1$.

When the arrival rate changes between periods or within periods, Green et al. (2001) proposed to adjust the value of the arrival rate λ in (12). This is a straightforward method that can give good staffing levels, at least for a call center operation that is similar to the $M(t)/M/s(t)$

Method	Λ_i
SIPPavg	$\int_{t_{i-1}}^{t_i} \lambda(t) dt$
SIPPmax	$\max_{t_{i-1} < t \leq t_i} \lambda(t)$
SIPPMix	If $\lambda(t)$ is nondecreasing in period i use SIPPavg rate, otherwise use SIPPmax rate.
LAGavg	$\int_{t_{i-1}-1/\mu}^{t_i-1/\mu} \lambda(t) dt$
LAGmax	$\max_{t_{i-1}-1/\mu < t \leq t_i-1/\mu} \lambda(t)$
LAGmix	If $\lambda(t)$ is nondecreasing in $[t_{i-1} - 1/\mu, t_i - 1/\mu]$ use LAGavg rate, otherwise use LAGmax rate.

Table 2: Different methods for adjusting the arrival rate to use in Equation (12). Here, t_i is the time when period i ends ($t_0 \equiv 0$), $1/\mu$ is the mean service time and Λ_i is the rate to be used to determine the staffing in period i .

model. They consider 6 different adjustments of the arrival rate. The 6 different schemes are given in Table 2. In Green et al. (2003) only SIPPavg, LAGavg and LAGmax are considered. When we computed the arrival rate to use for the LAG methods in the first period we assumed that the arrival rate prior to time zero was equal to the arrival rate at the beginning of the first period.

The required staffing, y_i , in period i is computed by letting $\lambda = \Lambda_i$ in (12) and letting

$$y_i = \min\{s > 0 \text{ and integer} : P(s) \geq 1 - 0.8\}.$$

The cost of the resulting staffing level is $f(y)$. The actual number of agents available in period i can actually be greater than y_i because of slack in the shift constraint $Ax \geq y$ in (1). We included the additional staffing from the slack when we evaluated the performance of the staffing levels obtained by these analytical heuristics.

5.4 Results

We enumerated the 16 experiments as in Table 3. The cost of the staffing levels obtained by each method is shown in Table 4.

Determining feasibility. The solutions obtained by the 6 queuing methods and the SAC-CPM are not guaranteed to be feasible for the call center staffing problem with expected service level constraints (2). It is true that a solution obtained by the SACCPM is feasible for the respective sample average approximation problem. To determine the feasibility we simulated each solution using a sample size of $n = 999$ (the maximum number of iterations in Promodel).

Since the simulation output of the 999 experiments is still random we decided to declare the service level requirements not met if the fraction of calls answered immediately is less than 75% in any period. This is similar to what Green et al. (2001) used to determine feasibility. They computed the service level using numerical integration of transient queuing formulas and said that the service level is infeasible in a period if the target probability is exceeded by at least 10%. We chose 75% as our threshold because the maximum 95% confidence interval half-width in our simulations using $n = 999$ was 4.4%.

The feasibility of the solutions is reported in Table 5. We first note that all the methods do well in most cases. The SIPPavg method struggles when there are no shifts present and the service rates are low. In this case there is a significant dependency between time periods which none of the SIPP methods take into account. We also note that the solutions from the SACCPM do not always meet the service level requirements by the 75% criterion and in many periods fail to meet the 80% requirements. This is because the sample size was only 100 when the solution was computed so there is a significant error in the estimates of the service level functions.

Ranking the methods. In Table 4 we put in boldface text the cost of the method that we declared a “winner.” We selected the winner based on two criteria. The first criterion is that the service level must be greater than 75% in all periods based on the simulation using sample size 999. The second criterion is cost. We see that the SACCPM is the winner in all but three experiments. In two of these the solutions fail to meet the feasibility criterion. In the third case the cost of the best solution was only 0.3% lower than the winner. In the case of ties it would make more sense to use one of the analytical queuing heuristics because they are much easier to implement. However, one would not know beforehand whether a tie would occur, and which heuristic to use.

The analytical queuing heuristics have been shown to do well on this particular type of queuing model. In the SACCPM, however, no assumptions are made on the arrival process or the service times, so it may apply in even more general settings.

In Section 4.1 we noted that the sample averages of the service level functions are not pseudoconcave for low staffing levels. In one or more iterations in some of the experiments we got zero pseudogradients, which cannot be used to generate feasibility cuts. Instead we imposed the lower bounds described in Section 4.1 and then checked upon termination of the algorithm whether these lower bounds were tight. The bounds were tight in experiment 16, so we relaxed the bounds and ran the algorithm until it terminated again, this time with a solution where these kinds of bounds were not tight.

5.5 Computational requirements

We can divide each iteration of the algorithm into 3 main parts in terms of computations:

1. *Solve the analytic center problem (7).* This usually took less than 1 second using the MINOS solver and never took more than 3 seconds.
2. *Solve the IP (8) to get an integer solution close to the analytic center.* In the beginning this takes on the order of seconds to solve. However, when there were no shifts, meaning that the solution space for y is larger, it could take millions of branch and bound nodes to find an optimal solution after a number of cuts were added. In fact it is not necessary to find an optimal solution of the IP (8) to advance the algorithm, so we often terminated the IP with a suboptimal, but feasible solution to determine the next iterate. If the IP seemed infeasible, but the infeasibility was difficult to prove and the optimality gap in the SACCPM was less than 1%, the SACCPM was terminated with a possibly suboptimal solution.
3. *Simulate to estimate the expected service level and gradients.* Simulating at each staffing level in ProModel took from 6 seconds to about 1 minute depending on the number of arrivals to the system, on a Pentium 4 3GHz computer. Up to 73 such simulations are required per iterations, although we did not always compute all 72 differences to compute the FD pseudogradient. It appeared that dependence between time periods was not a factor over more than 10 periods, so as a rule of thumb we only computed the differences for the 10 periods preceding period i and for period i if the service level constraint was not satisfied in period i . Staffing levels in subsequent periods do not have an effect on the service level in period i since the requirement is to answer 80% of the calls immediately.

Hence, estimating the pseudogradients and solving the IP (8) require the most computation; see Section 6 for ideas on how the computational requirements can be reduced.

The number of iterations required is given in Table 4. In the initial stages, the SACCPM sometimes produced several optimality cuts before any feasibility cuts were generated. Because the weight on the optimality cuts is only 1 in the beginning, the optimality cuts will be fairly close to each other in such a case and convergence will be slow. This happens when the starting point has much higher staffing levels than what is really needed. When this occurred, we added deeper optimality cuts until the first feasibility cuts were generated by the algorithm, i.e., we lowered z . One could start the algorithm close to the solutions of the analytical queuing heuristics, to hopefully speed up the convergence. We did not try to implement this approach.

Experiment	μ	R	θ	Shifts
1	4	8	0.75	Y
2	16	8	0.75	Y
3	4	32	0.75	Y
4	16	32	0.75	Y
5	4	8	0.25	Y
6	16	8	0.25	Y
7	4	32	0.25	Y
8	16	32	0.25	Y
9	4	8	0.75	N
10	16	8	0.75	N
11	4	32	0.75	N
12	16	32	0.75	N
13	4	8	0.25	N
14	16	8	0.25	N
15	4	32	0.25	N
16	16	32	0.25	N

Table 3: The experiments and the factor settings.

6 Conclusion and future research

The SACCPM presented in this paper is a promising method for computing optimal staffing levels in a call center when traditional methods fail. Although the computational requirements of the method are large, we were able to solve, or at least approximately solve, a number of moderately sized hypothetical call center staffing problems. The results of the SACCPM were encouraging and showed that the method can potentially be applied in real world situations with good results. Furthermore, the algorithm can be automated, so that no user intervention is required while the algorithm is running.

The algorithm is robust in the sense that it works well on a range of different problems which gives it further credibility over the traditional methods. Compared to the traditional queuing methods, the SACCPM does best in the presence of shift constraints when compared to the traditional queuing methods. To see why, recall that the SACCPM solves *both* the problem of determining minimum staffing levels *and* the problem of computing the best assignments to cover these staffing levels, while the queuing methods compute the required staffing levels first and the shift assignments afterwards, using the staffing levels as input.

There are several interesting directions for related future research, some of which we are actively pursuing. From the numerical experiments we identified that both the simulations and solving the integer programs can be computationally expensive. It would be beneficial to study

Experiment	SACCPM	SIPPavg	SIPPmax	SIPPMix	LAGavg	LAGmax	LAGmix	Iter
1	1008	1056 104.8%	1056 104.8%	1056 104.8%	1056 104.8%	1056 104.8%	1056 104.8%	17
2	1032	1056 102.3%	1056 102.3%	1056 102.3%	1032 100.0%	1056 102.3%	1032 100.0%	20
3	3456	3552 102.8%	3624 104.9%	3576 103.5%	3456 100.0%	3552 102.8%	3552 102.8%	10
4	3504	3552 101.4%	3624 103.4%	3576 102.1%	3576 102.1%	3576 102.1%	3528 100.7%	15
5	936	936 100.0%	936 100.0%	936 100.0%	936 100.0%	936 100.0%	936 100.0%	39
6	936	936 100.0%	936 100.0%	936 100.0%	936 100.0%	936 100.0%	936 100.0%	39
7	3024	3048 100.8%	3096 102.4%	3072 101.6%	3048 100.8%	3048 100.8%	3048 100.8%	12
8	2976	3048 102.4%	3096 104.0%	3072 103.2%	3024 101.6%	3072 103.2%	3048 102.4%	11
9	829	848 102.3%	862 104.0%	855 103.1%	848 102.3%	862 104.0%	855 103.1%	47
10	838	848 101.2%	858 102.4%	853 101.8%	847 101.1%	862 102.9%	853 101.8%	24
11	2746	2786 101.5%	2838 103.4%	2812 102.4%	2787 101.5%	2838 103.4%	2813 102.4%	38
12	2786	2786 100.0%	2838 101.9%	2812 100.9%	2778 99.7%	2830 101.6%	2804 100.6%	26
13	846	856 101.2%	862 101.9%	859 101.5%	856 101.2%	862 101.9%	859 101.5%	51
14	850	854 100.5%	860 101.2%	857 100.8%	856 100.7%	861 101.3%	859 101.1%	27
15	2774	2802 101.0%	2818 101.6%	2810 101.3%	2802 101.0%	2818 101.6%	2810 101.3%	32
16	2790	2802 100.4%	2818 101.0%	2810 100.7%	2797 100.3%	2815 100.9%	2806 100.6%	53

Table 4: Cost of the solutions and the number of iterations of the SACCPM. The bold numbers in each row are the lowest cost solutions out of the solutions that satisfy the feasibility requirements; see Table 5. The percentages are the percentage of the SACCPM cost in each experiment. The last column, “Iter”, shows the number of iterations of the SACCPM.

Experiment	SACCPM	SIPPavg	SIPPmax	SIPPMix	LAGavg	LAGmax	LAGmix
1	2 79.0%	1 79.0%	1 79.0%	2 79.0%	81.3%	84.4%	85.9%
2	82.1%	83.1%	83.1%	83.1%	82.1%	83.1%	82.1%
3	1 76.7%	4 1 74.4%	81.6%	81.5%	82.1%	85.9%	82.5%
4	82.1%	82.0%	84.3%	83.1%	82.0%	86.2%	82.1%
5	82.4%	83.1%	81.3%	82.5%	81.8%	83.3%	83.1%
6	81.8%	81.8%	81.8%	81.8%	81.8%	81.8%	81.8%
7	80.1%	80.4%	82.5%	80.8%	83.5%	83.8%	83.8%
8	3 76.6%	82.0%	82.5%	83.0%	81.6%	82.3%	81.6%
9	15 3 74.2%	14 3 73.2%	3 76.2%	3 76.2%	2 76.5%	81.2%	1 79.8%
10	12 77.8%	3 78.4%	80.8%	80.8%	2 79.9%	81.1%	2 79.9%
11	21 3 74.7%	33 23 61.2%	29 4 71.5%	29 4 71.5%	4 78.7%	80.8%	80.7%
12	5 77.9%	11 76.0%	80.8%	80.8%	1 79.7%	81.9%	80.2%
13	7 76.1%	1 79.9%	80.5%	80.5%	1 80.0%	1 80.0%	1 80.0%
14	6 75.9%	3 79.0%	1 79.8%	2 79.6%	2 79.6%	1 79.8%	2 79.6%
15	17 75.1%	19 76.7%	8 78.2%	8 78.2%	3 79.0%	80.3%	2 79.0%
16	10 76.6%	2 79.1%	81.1%	81.1%	1 79.7%	81.1%	81.1%

Table 5: Feasibility of the solutions. Each cell has up to 3 values. The top value is the number of periods in which the fraction of calls answered immediately is less than 80%. The center value is the number of periods in which the fraction of calls answered immediately is less than 75%. The bottom value shows the lowest fraction of calls answered immediately in any period. We do not display the top two values when they are equal to 0.

more efficient methods for estimating the pseudogradients that could mimic the performance of the finite difference method. In relation to the integer programs one should investigate integer programming algorithms that can utilize the special structure of the relaxed problems solved in each iteration and consider allowing approximate solutions of the IPs, especially in early stages of the algorithms. Another technique to speed up the computation of the next iterate in the continuous case is to drop cuts that are redundant (see, e.g., Mitchell, 2003, for more on this approach), and it is quite possible that the same technique could reduce the time required to solve the IPs in the later stages of the SACCPM. One could also generate some initial cuts by running simulations at the staffing levels suggested by heuristics.

It is quite possible that other optimization methods could perform well in this setting. The extended cutting plane method in Westerlund and Pörn (2002) seems to fit the framework particularly well, although some details of the implementation are unclear.

The problems solved in this paper were fairly simple instances of a call center staffing problem, but since no assumptions are made on the arrival and service processes and simulation is used to evaluate performance, it is quite possible that the method would also apply in more complicated settings. Call abandonments, skill-based routing and prioritizing multiple customer classes are problems that call center managers commonly face and it would be interesting to incorporate those in the algorithm. Moreover, it is likely that the SACCPM could, with appropriate modifications, be applied to problems other than call center staffing.

Acknowledgments

This research was supported by National Science Foundation grants DMI 0230528 and DMI 0400287 and a Horace H. Rackham School of Graduate Studies Faculty Grant. This research was conducted while the first author was a doctoral student at the Department of Industrial and Operations Engineering at the University of Michigan and he would like to thank the department for its generous financial support.

References

- S. Andradóttir. Simulation optimization. In J. Banks, editor, *Handbook of Simulation*, chapter 9, pages 307–333. John Wiley & Sons, New York, 1998.
- D. S. Atkinson and P. M. Vaidya. A cutting plane algorithm for convex programming that

- uses analytic centers. *Math. Programming*, 69(1, Ser. B):1–43, 1995. Nondifferentiable and large-scale optimization (Geneva, 1992).
- J. Atlason. *A Simulation Based Cutting Plane Method for Optimization of Service Systems*. PhD thesis, University of Michigan, Ann Arbor, MI, 2004.
- J. Atlason, M. A. Epelman, and S. G. Henderson. Using simulation to approximate subgradients of convex performance measures in service systems. In S. Chick, P. J. Sánchez, D. Ferrin, and D. J. Morrice, editors, *Proceedings of the 2003 Winter Simulation Conference*, pages 1824–1832, Piscataway, NJ, 2003. IEEE.
- J. Atlason, M. A. Epelman, and S. G. Henderson. Call center staffing with simulation and cutting plane methods. *Annals of Operations Research*, 127:333–358, 2004.
- O. Bahn, O. du Merle, J.-L. Goffin, and J.-P. Vial. A cutting plane method from analytic centers for stochastic programming. *Math. Programming*, 69(1, Ser. B):45–73, 1995. Nondifferentiable and large-scale optimization (Geneva, 1992).
- M. S. Bazaraa, H. D. Sherali, and C. M. Shetty. *Nonlinear Programming: Theory and Algorithms*. John Wiley & Sons, New York, 1993.
- S. Borst, A. Mandelbaum, and M. I. Reiman. Dimensioning large call centers. *Operations Research*, 52:17–34, 2004.
- I. Castillo, T. Joro, and Y. Li. Workforce scheduling with multiple objectives. Submitted, 2003.
- B. Cleveland and J. Mayben. *Call Center Management on Fast Forward*. Call Center Press, Annapolis, MD, 1997.
- O. duMerle. *Interior points and cutting planes: Development and implementation of methods for convex optimization and large scale structured linear programming. (In French)*. PhD thesis, University of Geneva, Geneva, Switzerland, 1995.
- O. duMerle, J.-L. Goffin, and J.-P. Vial. On improvements to the analytic center cutting plane method. *Computational Optimization and Applications*, 11:37–52, 1998.
- S. Elhedhli and J.-L. Goffin. The integration of an interior-point cutting plane method within a branch-and-price algorithm. *Math. Programming*, 100, Ser. A:267–294, 2003.
- M. C. Fu and J. Q. Hu. *Conditional Monte Carlo: Gradient Estimation and Optimization Applications*. Kluwer, Norwell, MA, 1997.

- N. Gans, G. Koole, and A. Mandelbaum. Telephone call centers: Tutorial, review and research prospects. *Manufacturing & Service Operations Management*, 5(2):79–141, 2003.
- P. Glasserman. *Gradient Estimation Via Perturbation Analysis*. Kluwer, Norwell, MA, 1991.
- P. W. Glynn. Likelihood ratio gradient estimation for stochastic systems. *Communications of the ACM*, 33:75–84, 1990.
- J.-L. Goffin, A. Haurie, and J.-P. Vial. Decomposition and nondifferentiable optimization with the projective algorithm. *Management Science*, 38(2):284–302, 1992.
- J.-L. Goffin, Z.-Q. Luo, and Y. Ye. Complexity analysis of an interior cutting plane method for convex feasibility problems. *SIAM J. Optim.*, 6(3):638–652, 1996.
- L. V. Green, P. J. Kolesar, and J. Soares. Improving the SIPP approach for staffing service systems that have cyclic demands. *Operations Research*, 49(4):549–564, 2001.
- L. V. Green, P. J. Kolesar, and J. Soares. An improved heuristic for staffing telephone call centers with limited operating hours. *Production and Operations Management*, 12(1):46–61, 2003.
- D. Gross and C. M. Harris. *Fundamentals of Queueing Theory*. John Wiley & Sons, New York, NY, 1998.
- A. Ingolfsson, E. Cabral, and X. Wu. Combining integer programming and the randomization method to schedule employees. Submitted, 2003.
- A. Ingolfsson, M. A. Haque, and A. Umnikov. Accounting for time-varying queueing effects in workforce scheduling. *European Journal of Operational Research*, 139:585–597, 2002.
- O. B. Jennings, A. Mandelbaum, W. A. Massey, and W. Whitt. Server staffing to meet time-varying demand. *Management Science*, 42(10):1383–1394, 1996.
- J.E. Kelley, Jr. The cutting-plane method for solving convex programs. *Journal of the Society for Industrial and Applied Mathematics*, 8(4):703–712, 1960.
- A. J. Kleywegt, A. Shapiro, and T. Homem-de-Mello. The sample average approximation method for stochastic discrete optimization. *SIAM Journal on Optimization*, 12(2):479–502, 2001.

- P. J. Kolesar and L. V. Green. Insights on service system design from a normal approximation to Erlang's delay formula. *Production and Operations Management*, 7:282–293, 1998.
- P. L'Ecuyer. A unified view on the IPA, SF, and LR gradient estimation techniques. *Management Science*, 36(11):1364–1383, 1990.
- P. L'Ecuyer. Note: On the interchange of derivative and expectation for likelihood ratio derivative estimators. *Management Science*, 41(4):738–748, 1995.
- A. Mandelbaum. Call centers (centres): Research bibliography with abstracts. Version 5. Accessible online via <http://ie.technion.ac.il/serveng/References/ccbib.pdf> [accessed June 7, 2004], 2003.
- A. J. Mason, D. M. Ryan, and D. M. Panton. Integrated simulation, heuristic and optimisation approaches to staff scheduling. *Operations Research*, 46:161–175, 1998.
- J. E. Mitchell. Computational experience with an interior point cutting plane algorithm. *SIAM J. Optim.*, 10(4):1212–1227 (electronic), 2000.
- J. E. Mitchell. Polynomial interior point cutting plane methods. *Optim. Methods Softw.*, 18(5):507–534, 2003.
- K. Murota. *Discrete Convex Analysis*. SIAM, Philadelphia, PA, 2003.
- Y. Nesterov. Complexity estimates of some cutting plane methods based on the analytic barrier. *Math. Programming*, 69(1, Ser. B):149–176, 1995. Nondifferentiable and large-scale optimization (Geneva, 1992).
- O. Peton and J.-P. Vial. A tutorial on ACCPM: User's guide for version 2.01. Working Paper. University of Geneva, 2001.
- R. Y. Rubenstein and A. Shapiro. *Discrete Event Systems: Sensitivity Analysis and Stochastic Optimization by the Score Function Method*. John Wiley & Sons, Chichester, England, 1993.
- A. Shapiro. Monte Carlo sampling methods. In A. Ruszczyński and A. Shapiro, editors, *Stochastic Programming*, Handbooks in Operations Research and Management Science. Elsevier, 2003. To appear.
- T. Westerlund and R. Pörn. Solving pseudo-convex mixed integer optimization problems by cutting plane techniques. *Optimization and Engineering*, 3:253–280, 2002.

W. Whitt. The pointwise stationary approximation for $M_t/M_t/s$ queues is asymptotically correct as the rates increase. *Management Science*, 37:307–314, 1991.