

# MONOTONE OPTIMAL POLICIES FOR A TRANSIENT QUEUEING STAFFING PROBLEM

MICHAEL C. FU

*Robert H. Smith School of Business, University of Maryland, College Park, Maryland 20742, mfu@rhsmith.umd.edu*

STEVEN I. MARCUS

*Department of Electrical Engineering, University of Maryland, College Park, Maryland 20742, marcus@isr.umd.edu*

I-JENG WANG

*The Johns Hopkins University Applied Physics Laboratory, Laurel, Maryland 20723,  
I-Jeng.Wang@jhuapl.edu*

(Received April 1997; revisions received March 1998, April 1998; accepted June 1998)

We consider the problem of determining the optimal policy for staffing a queueing system over multiple periods, using a model that takes into account transient queueing effects. Formulating the problem in a dynamic programming setting, we show that the optimal policy follows a monotone optimal control by establishing the submodularity of the objective function with respect to the staffing level and initial queue size in a period. In particular, this requires proving that the system occupancy in a  $G/M/s$  queue is submodular in the number of servers and initial system occupancy.

In Yoo (1996), a dynamic programming (DP) model is formulated to address the problem of setting staffing levels at a post office's service windows over multiple time periods (typically of length 15 minutes or 30 minutes) in a day. A major component of the model is the inclusion of transient queueing behavior—obtained via numerical integration—so that deviation from traditional steady-state models could be investigated. A computationally efficient heuristic algorithm is proposed to solve the DP. One of the key assumptions that underlies the algorithm is monotonicity of the optimal policy, which results in a significant reduction in searching the state space. The purpose of this note is to establish this structural property. The chief result needed to accomplish this, which is interesting in its own right, is the submodularity of the system occupancy at any point in time with respect to the initial occupancy and the number of servers.

To be specific, we consider a single-queue system with multiple servers, in which the arrival process of customers is nonstationary. The objective is to choose staffing levels over a multiperiod time horizon (during each period, the arrival rate is assumed approximately constant) to minimize the expected cost comprising a weighted sum of server costs and queue costs, where queue costs are evaluated using transient queueing models. Aside from making the performance analysis more difficult, the inclusion of transient queueing effects also complicates the staffing problem. Periods are no longer uncoupled because (a) there is a dependence on initial conditions, and (b) the system does not necessarily reach steady state by the end of a period; hence, staffing decisions for different periods cannot be implemented independently. The experience with the post office setting reported in Yoo (1996)

is that the manager sets discrete intervals at which to make staffing decisions, but he/she also retains the flexibility to handle unexpected extreme surges or lulls by adding/subtracting staff at any time. Our model primarily addresses the former decision-making problem. Thus we consider an operational version of the problem in which a decision is made at the beginning of each period as to how many servers to staff for that period, given the number of customers in the system at the beginning of the period. Specifically, the problem is formulated as a finite-horizon dynamic program.

When the goal is numerical solutions of the DP, structural properties can lead to great gains in computational efficiency. Monotone optimal policies constitute one of the most well-known and useful of such characterizations; see, for example, Heyman and Sobel (1984, Chapter 8) or Puterman (1994, §4.7.3). In terms of the staffing problem, this monotonicity translates into the intuitively obvious result that the optimal policy should prescribe increased staffing levels for higher initial queue lengths.

The main property needed to establish a monotone optimal policy is submodularity of the period cost function (Heyman and Sobel 1984). To establish submodularity for our problem requires proving that the system occupancy is submodular in the number of servers and initial system occupancy. We prove this result for a  $G/M/s$  queue. A closely related result is given in Chang et al. (1994, theorem 5.3.21), in which it is shown that the queue length process for a pure death process is stochastically increasing and supermodular in the initial state and time. Our result here establishes an analogous submodularity result with respect to the initial number in the

*Subject classifications:* Dynamic programming, applications: staffing problem. Queues, transient results: submodularity. Optimal control: monotone policies.  
*Area of review:* STOCHASTIC MODELS.

system and the number of servers, for any fixed time horizon. The proof uses uniformization, as in Chang et al. (1994).

## 1. THE DYNAMIC PROGRAMMING FORMULATION

The state of the system is the number of customers in the system—in service and in queue—referred to as the *system occupancy* throughout this paper. The action taken each period is the number of servers to staff for that entire period and is based only on the state of the system at the beginning of each period. In sum, our model is discrete time and discrete state space. The period cost is the sum of the staffing cost and the expected time-averaged system occupancy for the period. The objective is to find the optimal policy that minimizes the total horizon costs.

We introduce the following notation to define the problem:

- $N$  Number of periods in the total horizon,
- $\tau$  Length of a period,
- $S_n$  Number of servers used in period  $n$ ,
- $X_n$  System state (occupancy) at the beginning of period  $n$ ,
- $s_{\max}$  Maximum number of servers that can be used in a period,
- $\mu$  Service rate for each server,
- $k$  Cost per server per period used,
- $\bar{c}_n(i, s)$  Cost function for period  $n$ , given  $i$  initial customers and  $s$  servers,
- $\sigma_n(i)$  Number of servers to be assigned to period  $n$ , if the initial system occupancy is  $i$  customers,
- $\sigma_n$  Vector of  $\sigma_n(i)$ , staffing policy for period  $n$ , and
- $\sigma$  Vector of  $\sigma_n$ , matrix denoting staffing policy for entire horizon.

To reduce notation, we have assumed that the length of each period is identical. The optimization problem can be stated as

$$\min_{\sigma} \sum_{n=1}^N E[\bar{c}_n(X_n, \sigma_n(X_n))], \quad (1)$$

where the distribution of  $X_n$  will depend on the policy  $\sigma$ . In order to characterize this quantity, we introduce the random variable (stochastic process)  $X_t^{(n)}(i, s)$  = system occupancy at time  $t$  into period  $n$ , given a system occupancy of  $i$  at the beginning of period  $n$  and  $s$  servers throughout. The period cost function is the sum of the expected time-averaged mean system occupancy and a linear cost on the number of servers:

$$\bar{c}_n(i, s) = \frac{1}{\tau} \int_0^{\tau} E[X_t^{(n)}(i, s)] dt + ks. \quad (2)$$

When we are interested in the system occupancy at the end of a period, we use the simplified notation

$$X^{(n)}(i, s) = X_{\tau}^{(n)}(i, s), \quad (3)$$

i.e., we drop the  $\tau$  subscript. In this case, we have  $X_n = X^{(n-1)}(X_{n-1}, S_{n-1})$ .

The solution to the adaptive staffing problem (1) can be obtained by applying stochastic DP. The state of a period is the number of customers at the beginning of the current period, which is equal to the number at the end of the previous period. The transition probabilities from period to period are governed by the transient behavior of the associated multiserver queue. We index periods in chronological time (i.e., forward), so the backward dynamic programming algorithm begins with the last period. The resulting DP recursive equation (optimality equation) for (1) is given by the following ( $n = 1, \dots, N$ ):

$$f_n^*(i) = \min_s \{J_n(i, s)\}, \quad (4)$$

$$J_n(i, s) = \bar{c}_n(i, s) + E[f_{n+1}^*(X^{(n)}(i, s))], \quad (5)$$

where we assume the final condition  $f_{N+1}^*(i) = 0$  for all  $i$ . The optimal policy is computed by backward induction on  $n = N, N-1, \dots, 1$ :

$$\sigma_n^*(i) = \arg \min_s \{\bar{c}_n(i, s) + E[f_{n+1}^*(X^{(n)}(i, s))]\}. \quad (6)$$

We now state our main result, which we will prove in the remainder of the paper.

**THEOREM 1.** *Under the assumption of exponential service times and interarrival times, there exists an optimal policy that is monotonically increasing in the state, i.e., such that for all  $i$  and  $n$ ,*

$$\sigma_n^*(i) \leq \sigma_n^*(i+1).$$

The monotone structure of the policy allows the search space to be reduced considerably, serving as the basis for the efficient heuristic algorithm proposed in Yoo (1996). Specifically, when (6) is evaluated to solve for the optimal policy, the computationally expensive portion is the evaluation of  $\bar{c}_n(i, s)$  for each  $s$ , given by (2), via Runge-Kutta numerical integration. If  $M$  is the maximum number of servers, then the result allows the search at state  $i$  to proceed from  $s = \sigma_n^*(i), \dots, M$ , instead of from  $s = 1, \dots, M$ . Furthermore, the algorithm usually requires only two evaluations at  $s = \sigma_n^*(i), \sigma_n^*(i) + 1$ , for the parameter settings typical in the post office setting, leading to an order-of-magnitude computational savings. Thus, problems for an 8-hour day of 15-minute intervals (giving 32 periods), which otherwise would have taken several hours to solve, are reduced to a few minutes of computation time. Alternatively, the model could also be used to approximate a continuous-time model over a much shorter horizon where the arrival process is stationary (e.g., by taking 30 periods over a 15-minute horizon, giving 30-second intervals), which might be well approximated by an infinite-horizon model with a stationary optimal policy.

Before proceeding, we note that the submodularity result proven in the next section actually holds for any general arrival process independent of the service times.

**2. MONOTONE OPTIMAL POLICIES**

Throughout, we use the monotonicity terms *increasing* and *decreasing* in the nonstrict sense. We first review the definitions of stochastic ordering and submodularity (cf. Shaked and Shanthikumar 1994, Ross 1983, and Stoyan 1983).

DEFINITION. A random variable  $X$  is stochastically smaller than a random variable  $Y$ , written  $X \leq_{st} Y$ , if

$$P(X > x) \leq P(Y > x) \quad \text{for all } x.$$

Two random variables are equal in distribution, written  $X =_{st} Y$ , if and only if  $X \leq_{st} Y$  and  $X \geq_{st} Y$ .

DEFINITION. A function  $\phi: \Re^2 \rightarrow \Re$  is *submodular* if for any  $x_1 < x_2, y_1 < y_2$ ,

$$\phi(x_1, y_1) + \phi(x_2, y_2) \leq \phi(x_1, y_2) + \phi(x_2, y_1).$$

We first define the stochastic version of submodularity in the natural way.

DEFINITION. A family of random variables  $\{X(i, s)\}$  is *stochastically submodular* if for any  $i_1 < i_2, s_1 < s_2$  there exist on a common probability space  $(\Omega, F, P)$  four random variables  $\hat{X}^j, j = 1, 2, 3, 4$ , equal in distribution to  $X(i_1, s_1), X(i_2, s_1), X(i_1, s_2), X(i_2, s_2)$ , respectively, such that for all  $\omega \in \Omega$ ,

$$\hat{X}^1(\omega) + \hat{X}^4(\omega) \leq \hat{X}^2(\omega) + \hat{X}^3(\omega). \tag{7}$$

In particular, the definition implies that the expectation is submodular in the ordinary sense.

For our application, we will need a stronger version.

DEFINITION. A family of random variables  $\{X(i, s)\}$  is *strongly stochastically submodular* if for any  $i_1 < i_2, s_1 < s_2$  there exist on a common probability space  $(\Omega, F, P)$  four random variables  $\hat{X}^j, j = 1, 2, 3, 4$ , equal in distribution to  $X(i_1, s_1), X(i_2, s_1), X(i_1, s_2), X(i_2, s_2)$ , respectively, such that for all  $\omega \in \Omega$ ,

$$\hat{X}^3(\omega) \leq \min\{\hat{X}^1(\omega), \hat{X}^4(\omega)\},$$

$$\hat{X}^1(\omega) + \hat{X}^4(\omega) \leq \hat{X}^2(\omega) + \hat{X}^3(\omega).$$

The strong version is analogous to the definition of stochastic submodularity given in Chang et al. (1994, definition 5.3.12), in that stochastic monotonicity in both arguments is also specified (and must hold on the *same* probability space). We require stochastic increasing in the first variable and stochastic decreasing in the second variable. Clearly, this strong stochastic version includes the previous version. We now present the following result, which will be needed later.

LEMMA 1. *If  $X(\cdot, \cdot)$  is strongly stochastically submodular, then  $P(X(\cdot, \cdot) > x)$  is submodular for all  $x$ .*

PROOF. By definition of strongly stochastic submodular, we have the existence of four random variables  $\hat{X}^j, j = 1, 2, 3, 4$ , defined on a common probability space  $(\Omega, F, P)$  and equal in distribution to  $X(i_1, s_1), X(i_2, s_1), X(i_1, s_2), X(i_2, s_2)$ , respectively,  $i_1 < i_2, s_1 < s_2$ , such that for all  $\omega \in \Omega$ ,

$$\hat{X}^3(\omega) \leq \min\{\hat{X}^1(\omega), \hat{X}^4(\omega)\},$$

$$\hat{X}^1(\omega) + \hat{X}^4(\omega) \leq \hat{X}^2(\omega) + \hat{X}^3(\omega).$$

Consider the following construction:

$$\hat{Y}^j(x) = \mathbf{1}\{\hat{X}^j > x\}, \quad j = 1, 2, 3,$$

$$\hat{Y}^4(x) = \mathbf{1}\{\hat{X}^3 > x\} + \mathbf{1}\{\hat{X}^1 < x \leq \hat{X}^1 + \hat{X}^4 - \hat{X}^3\}$$

$$=_{st} \mathbf{1}\{\hat{X}^4 > x\},$$

because  $\hat{X}^3 \leq \hat{X}^1$  and  $\hat{X}^3 \leq \hat{X}^4$ , where  $\mathbf{1}\{\cdot\}$  is the set indicator function. Then it can be easily verified that  $\hat{Y}^j(x)$  satisfies the required inequality for stochastic submodularity; i.e., for any  $x$ , we have for all  $\omega$ ,

$$\hat{Y}^1(x) + \hat{Y}^4(x) \leq \hat{Y}^2(x) + \hat{Y}^3(x),$$

so that taking the expectation yields the result:

$$P(\hat{X}^1 > x) + P(\hat{X}^4 > x) \leq P(\hat{X}^2 > x) + P(\hat{X}^3 > x). \quad \square$$

We use the following, fairly standard characterization of the existence of monotone optimal policies to establish our main result.

PROPOSITION 1 (Theorem 8-5 in Heyman and Sobel 1984). *Suppose the following hold for the dynamic programming problem defined by (1)–(5):*

- (i)  $\bar{c}_n(\cdot, s)$  is increasing for all  $s$  and  $n$ ;
- (ii)  $\bar{c}_n(\cdot, \cdot)$  is submodular and bounded below for all  $n$ ;
- (iii)  $P(X^{(n)}(\cdot, \cdot) > x)$  is submodular for all  $x$  and  $n$ ;
- (iv)  $X^{(n)}(\cdot, s)$  is stochastically increasing for all  $s$  and  $n$ .

*Then, for each  $n$ , there exists  $\sigma_n^*(\cdot)$  increasing that is optimal, i.e., satisfies (6).*

The statement of the result has been modified slightly to take into account the fact that for our setting, the action space  $\{1, \dots, s_{\max}\}$  is finite and independent of the state, so certain conditions are automatically satisfied. In particular, the action space is compact, contracting, and ascending; the minimum is achieved in (4); and  $J_n(i, \cdot)$  is lower semicontinuous.

A similar result is given by Puterman (1994, theorem 4.7.4) using the terminology subadditive/superadditive instead of submodular/supermodular (and considering maximization instead of minimization).

In order to apply this result to prove our main theorem, we make the following observations:

- $\bar{c}_n(\cdot, \cdot)$  is bounded below by 0;
- Integration preserves stochastic ordering and also submodularity;
- $\bar{c}_n(\cdot, \cdot)$  is composed of a linear cost and an expectation of a time-averaged integral of  $X^{(n)}(\cdot, \cdot)$ ;
- Conditions (iii) and (iv) are implied by strong stochastic submodularity, via the definition and Lemma 1.

We note that the stochastic increasing property with respect to initial system occupancy is also proven for the general case (no exponential assumptions) using a sample path proof in Assad et al. (1997); see also Shanthikumar

and Yao (1989) and Sonderman (1979). We now proceed to establish that  $X(\cdot, \cdot)$  is strongly stochastically submodular, where we drop the period-dependent superscript for notational brevity. We require exponential service times, but the arrival process can be general, as long as it is independent of the service times.

LEMMA 2. *The system occupancy at any time in a G/M/s queue is strongly stochastically submodular with respect to the initial system occupancy and the number of servers, so that the tail distribution of system occupancy is submodular.*

PROOF. We begin with a pure death process and then superimpose the arrival process. Because of the memoryless property of the exponentially distributed service times, this presents little additional difficulty. The pure death process is handled by uniformizing the process at rate  $(s + 1)\mu$  and defining four processes on the same probability space.

Without loss of generality, we take  $\mu = 1$  throughout. Because the process is constant between uniformized epochs, it suffices to establish the result for the discrete uniformized epochs, which we denote by  $t_0 = 0, t_1, t_2, \dots$ . We proceed by induction on  $m$  for epoch  $t_m$ . Because  $X(\cdot, \cdot)$  is defined only on the integers, it suffices to establish the requisite relationships for  $i_1 = i, i_2 = i + 1, s_1 = s, s_2 = s + 1$ . For convenience, we simplify the notation a bit by introducing the following definitions:

$$\begin{aligned} Y_m^1 &= X_{t_m}(i, s), \\ Y_m^2 &= X_{t_m}(i + 1, s), \\ Y_m^3 &= X_{t_m}(i, s + 1), \\ Y_m^4 &= X_{t_m}(i + 1, s + 1). \end{aligned}$$

We will, in particular, show the somewhat stronger result—that there exist on a common probability space  $(\Omega, F, P)$  four random variables  $\hat{Y}_m^j, j = 1, 2, 3, 4$ , equal in distribution to  $Y_m^j, j = 1, 2, 3, 4$ , respectively, such that for all  $\omega \in \Omega$  and for all  $m$ ,

$$\begin{aligned} \hat{Y}_m^1(\omega) &\leq \hat{Y}_m^2(\omega), & \hat{Y}_m^3(\omega) &\leq \hat{Y}_m^4(\omega); \\ \hat{Y}_m^3(\omega) &\leq \hat{Y}_m^1(\omega); \\ \hat{Y}_m^1(\omega) + \hat{Y}_m^4(\omega) &\leq \hat{Y}_m^2(\omega) + \hat{Y}_m^3(\omega). \end{aligned}$$

Clearly, for  $m = 0$ , the result holds because  $X_{t_0}(i, s) = i$  for any  $s$ ; hence,

$$\begin{aligned} Y_0^1 &= X_{t_0}(i, s) = i < i + 1 = X_{t_0}(i + 1, s) = Y_0^2, \\ Y_0^3 &= X_{t_0}(i, s + 1) = i < i + 1 = X_{t_0}(i + 1, s + 1) = Y_0^4, \\ Y_0^3 &= X_{t_0}(i, s + 1) = i = X_{t_0}(i, s) = Y_0^1, \\ Y_0^1 + Y_0^4 &= X_{t_0}(i, s) + X_{t_0}(i + 1, s + 1) = 2i + 1 \\ &= X_{t_0}(i + 1, s) + X_{t_0}(i, s + 1) = Y_0^2 + Y_0^3. \end{aligned}$$

Now assuming the result holds for  $m$ , we establish it for  $m + 1$ .

The four random variables at  $m + 1$  are related to those at  $m$  via as follows:

$$Y_{m+1}^j = Y_m^j - C_m^j, \quad j = 1, 2, 3, 4,$$

where  $C_m^j$  are Bernoulli random variables, i.e., 1 w.p.  $p_m^j$  and 0 otherwise, where

$$p_m^1 = (s \wedge Y_m^1)/(s + 1), \quad (8)$$

$$p_m^2 = (s \wedge Y_m^2)/(s + 1), \quad (9)$$

$$p_m^3 = ((s + 1) \wedge Y_m^3)/(s + 1), \quad (10)$$

$$p_m^4 = ((s + 1) \wedge Y_m^4)/(s + 1), \quad (11)$$

where  $\wedge$  denotes the minimum function, i.e.,  $x \wedge y = \min(x, y)$ .

By the induction hypothesis on  $m$ , there exist on a common probability space  $(\Omega, F, P)$  four random variables  $\hat{Y}_m^j, j = 1, 2, 3, 4$ , equal in distribution to  $Y_m^j, j = 1, 2, 3, 4$ , respectively, such that for all  $\omega \in \Omega$ ,

$$\hat{Y}_m^1(\omega) \leq \hat{Y}_m^2(\omega), \quad \hat{Y}_m^3(\omega) \leq \hat{Y}_m^4(\omega); \quad (12)$$

$$\hat{Y}_m^3(\omega) \leq \hat{Y}_m^1(\omega); \quad (13)$$

$$\hat{Y}_m^1(\omega) + \hat{Y}_m^4(\omega) \leq \hat{Y}_m^2(\omega) + \hat{Y}_m^3(\omega). \quad (14)$$

Let  $U \sim U(0, s + 1)$  be a random variable defined on the same probability space as  $\hat{Y}_m^j$ ; i.e.,  $U$  will be the uniformization random variable uniformly distributed on  $[0, s + 1]$  (recall  $\mu = 1$ ). Define four random variables  $\hat{C}_m^j, j = 1, \dots, 4$ , on this same probability space by

$$\hat{C}_m^j = \mathbf{1}\{U \leq s \wedge \hat{Y}_m^j\}, \quad j = 1, 2, \quad (15)$$

$$\hat{C}_m^3 = \mathbf{1}\{U \leq (s + 1) \wedge \hat{Y}_m^3\}, \quad (16)$$

$$\begin{aligned} \hat{C}_m^4 &= \mathbf{1}\{U \leq (s + 1) \wedge \hat{Y}_m^3\} \\ &\quad + \mathbf{1}\{(s + 1) \wedge \hat{Y}_m^3 < U \leq s \wedge \hat{Y}_m^1 + (s + 1) \wedge \hat{Y}_m^4 - (s + 1)\} \\ &\quad + \mathbf{1}\{s \wedge \hat{Y}_m^1 < U \leq s \wedge \hat{Y}_m^1 + (s + 1) \wedge \hat{Y}_m^4 \\ &\quad - (s + 1) \wedge \hat{Y}_m^3\}. \end{aligned} \quad (17)$$

Note that for  $\hat{C}_m^4$ , the three indicator functions are defined on mutually exclusive events unless  $\hat{Y}_m^3(\omega) \geq s + 1$ , but in this case

$$\begin{aligned} \hat{Y}_m^3(\omega) \geq s + 1 &\Rightarrow \hat{Y}_m^4(\omega) \geq s + 1 \\ &\Rightarrow (s + 1) \wedge \hat{Y}_m^4(\omega) = (s + 1) \wedge \hat{Y}_m^3(\omega) = s + 1, \end{aligned}$$

so that both the second and third indicators would be zero. Note also that the second indicator would be zero whenever  $s \wedge \hat{Y}_m^1 + (s + 1) \wedge \hat{Y}_m^4 - (s + 1) \wedge \hat{Y}_m^3 \leq s + 1$ . Thus, we have  $\hat{C}_m^j =_{\text{st}} C_m^j$  for  $j = 1, \dots, 4$ , and we take the construction

$$\hat{Y}_{m+1}^j = \hat{Y}_m^j - \hat{C}_m^j,$$

where  $\hat{Y}_{m+1}^j =_{\text{st}} Y_{m+1}^j$ . For ease of presentation, we use  $(y_1, \dots, y_4, c_1, \dots, c_4, u)$  to denote a sample of  $(\hat{Y}_m^1, \dots, \hat{Y}_m^4, \hat{C}_m^1, \dots, \hat{C}_m^4, U)$ , and  $(y'_1, \dots, y'_4)$  a sample of  $(\hat{Y}_{m+1}^1, \dots, \hat{Y}_{m+1}^4)$ , so  $y'_j = y_j - c_j, j = 1, 2, 3, 4$ , where  $c_j$  depends on  $u$  and  $y_j$  (and also  $y_1$  and  $y_3$  for  $j = 4$ ). We now establish

$$y'_1 \leq y'_2, \quad y'_3 \leq y'_4; \quad (18)$$

$$y'_3 \leq y'_1; \quad (19)$$

$$y'_1 + y'_4 \leq y'_2 + y'_3. \quad (20)$$

We first establish the inequalities stated in (18) and (19). Note that because  $c_i - c_j \leq 1$  for all  $i, j = 1, \dots, 4$ , it suffices to show each inequality of  $y'_i$  and  $y'_j$  in (18) or (19) assuming that  $y_i = y_j$ . In this case, we have  $c_1 = c_2$  and  $c_3 = c_4$ . Thus  $y'_1 = y'_2$  and  $y'_3 = y'_4$ . To show (19), we observe that  $c_3 \geq c_1$  because  $(s+1) \wedge y_3 \geq s \wedge y_1$  when  $y_1 = y_3$ . Therefore, we have  $y'_3 \leq y'_1$ .

Now we show that inequality (20) holds. First, (20) translates into

$$y_1 + y_4 - (c_1 + c_4) \leq y_2 + y_3 - (c_2 + c_3). \quad (21)$$

Because by induction we have  $y_1 + y_4 \leq y_2 + y_3$ , we consider two possibilities, depending on whether the strict inequality holds or equality holds. If  $y_1 + y_4 < y_2 + y_3$ , then we have  $y_1 + y_4 \leq y_2 + y_3 - 1$ , and the only way (21) could not hold is if  $c_2 = c_3 = 1$  (recall  $c_i$  is either 0 or 1). By (12) and (13) and the construction (15)–(17), we have  $c_3 \leq c_4$ , so  $c_3 = 1 \Rightarrow c_4 = 1 \Rightarrow c_1 + c_4 \geq 1$ , and hence (21) holds.

Thus it remains to be shown that

$$c_1 + c_4 \geq c_2 + c_3 \quad (22)$$

for the case where

$$y_1 + y_4 = y_2 + y_3. \quad (23)$$

By the construction of  $\hat{C}_m^j$  given by (15)–(17), we have

$$\begin{aligned} c_2 - c_1 &= \mathbf{1}\{s \wedge y_1 < u \leq s \wedge y_2\}, \\ c_4 - c_3 &= \mathbf{1}\{s \wedge y_1 < u \leq s \wedge y_1 + (s+1) \wedge y_4 \\ &\quad - (s+1) \wedge y_3\} \\ &\quad + \mathbf{1}\{(s+1) \wedge y_3 < u \leq s \wedge y_1 + (s+1) \wedge y_4 \\ &\quad - (s+1)\}, \end{aligned}$$

so (22) holds if

$$s \wedge y_1 + (s+1) \wedge y_4 \geq s \wedge y_2 + (s+1) \wedge y_3. \quad (24)$$

We pursue the proof by considering the following cases for each sample (recall  $y_1 \leq y_2$ ):

(i)  $s \leq y_1 \leq y_2$ : Hence  $s \wedge y_2 = s \wedge y_1 = s$ . By (12) we also have  $y_4 \geq y_3$ , so  $(s+1) \wedge y_4 \geq (s+1) \wedge y_3$ . Therefore inequality (24) holds.

(ii)  $s \geq y_2 \geq y_1$ : Hence  $s \wedge y_1 = (s+1) \wedge y_1$ , and  $s \wedge y_2 = (s+1) \wedge y_2$ .

Define a function  $\zeta : \mathfrak{R} \rightarrow \mathfrak{R}$  by

$$\zeta(x) = (s+1) \wedge x.$$

Because  $\zeta(\cdot)$  is concave and increasing, we have by (23)

$$\zeta(y_2) + \zeta(y_3) \leq \zeta(y_1) + \zeta(y_4).$$

Therefore, inequality (24) holds.

(iii)  $y_2 > s > y_1$ : Hence  $s \wedge y_1 = y_1$ , and  $s \wedge y_2 = s$ . By (13) we also have  $s > y_1 \geq y_3$ , so  $(s+1) \wedge y_3 = y_3$ .

We further consider the following two cases:

(a)  $y_4 \geq s+1$ : Hence  $(s+1) \wedge y_4 = s+1$ , so

$$\begin{aligned} s \wedge y_1 + (s+1) \wedge y_4 &= y_1 + (s+1) \\ &> y_3 + s \\ &= s \wedge y_2 + (s+1) \wedge y_3. \end{aligned}$$

(b)  $y_4 < s+1$ : Hence  $(s+1) \wedge y_4 = y_4$ , so by (23), we have

$$\begin{aligned} s \wedge y_1 + (s+1) \wedge y_4 &= y_1 + y_4 \\ &= y_2 + y_3 \\ &> s + y_3 \\ &= s \wedge y_2 + (s+1) \wedge y_3. \end{aligned}$$

By (i)–(iii), we have established (20), completing the proof of strong stochastic submodularity for the pure death process. The progression of the proof is such that superimposing an arrival process does not change the main line of proof based on induction on event (arrival or departure) epochs. At arrival epochs, the relevant relationships will still hold because 1 is added to both sides; and at departure epochs, the same argument can be used as in the pure death process.  $\square$

Combining Proposition 1 with Lemmas 1 and 2 finishes the proof of the theorem.

## ACKNOWLEDGMENT

This work was supported in part by the National Science Foundation under Grant No. NSF EEC 94-02384. The authors thank David Yao for suggesting the construction used in the stochastic submodularity proof of Lemma 1 and the form of construction used in the proof of Lemma 2.

## REFERENCES

- Assad, A. A., M. C. Fu, J. S. Yoo. 1997. A lower bounding result for the optimal policy in an adaptive staffing problem. ISR Technical Report, University of Maryland, College Park, MD.
- Chang, C. S., J. G. Shanthikumar, D. D. Yao. 1994. Stochastic convexity and stochastic majorization. *Stochastic Modeling and Analysis of Manufacturing Systems*. David D. Yao, ed. Springer-Verlag, New York.
- Heyman, S. P., M. J. Sobel. 1984. *Stochastic Models in Operations Research*, Volume II. McGraw-Hill, New York.
- Puterman, M. L. 1994. *Markov Decision Processes*. John Wiley & Sons, New York.
- Ross, S. 1983. *Stochastic Processes*. Wiley, New York.
- Shaked, M., G. Shanthikumar. 1994. *Stochastic Orders and Their Applications*. Academic Press, Boston, MA.
- Shanthikumar, J., D. Yao. 1989. Stochastic monotonicity in general queueing networks. *J. Appl. Probab.* **26** 413–417.
- Sonderman, D. 1979. Comparing multi-server queues with finite waiting rooms—II: different number of servers. *Adv. in Appl. Probab.* **11** 448–455.
- Stoyan, D. 1983. *Comparison Methods for Queues and Other Stochastic Models*. John Wiley, New York.
- Yoo, J. 1996. *Queueing models for staffing service operations*. Ph.D. dissertation. University of Maryland, College Park, MD.