



## Management Science

Publication details, including instructions for authors and subscription information:  
<http://pubsonline.informs.org>

### Note—Optimal Assignment of Customers in a Two Server Congestion System with No Waiting Room

Wayne Winston,

To cite this article:

Wayne Winston, (1978) Note—Optimal Assignment of Customers in a Two Server Congestion System with No Waiting Room. Management Science 24(6):702-705. <https://doi.org/10.1287/mnsc.24.6.702>

Full terms and conditions of use: <https://pubsonline.informs.org/Publications/Librarians-Portal/PubsOnLine-Terms-and-Conditions>

This article may be used only for the purposes of research, teaching, and/or private study. Commercial use or systematic downloading (by robots or other automatic processes) is prohibited without explicit Publisher approval, unless otherwise noted. For more information, contact [permissions@informs.org](mailto:permissions@informs.org).

The Publisher does not warrant or guarantee the article's accuracy, completeness, merchantability, fitness for a particular purpose, or non-infringement. Descriptions of, or references to, products or publications, or inclusion of an advertisement in this article, neither constitutes nor implies a guarantee, endorsement, or support of claims made of that product, publication, or service.

© 1978 INFORMS

Please scroll down for article—it is on subsequent pages



With 12,500 members from nearly 90 countries, INFORMS is the largest international association of operations research (O.R.) and analytics professionals and students. INFORMS provides unique networking and learning opportunities for individual professionals, and organizations of all types and sizes, to better understand and use O.R. and analytics tools and methods to transform strategic visions and achieve better outcomes.

For more information on INFORMS, its publications, membership, or meetings visit <http://www.informs.org>



## OPTIMAL ASSIGNMENT OF CUSTOMERS IN A TWO SERVER CONGESTION SYSTEM WITH NO WAITING ROOM†

WAYNE WINSTON‡

We consider a two server congestion system which is heterogeneous in the sense that the reward received depends on the match between server and customer. In particular, if a type  $s$  customer ( $s = 1, 2, 3, \dots$ ) is assigned to server  $i$  ( $i = 1, 2$ ) a reward  $R_{si}$  is earned. Service time is assumed to depend only on customer type. Upon arrival a customer must be assigned to a server, (if one is available) with all customers who find both servers occupied being turned away. The policy that maximizes the long-run expected reward earned over an infinite horizon is shown to depend on a single critical number. Applications to the deployment of fire engines and assignment of patients to coronary care units are briefly discussed.

### 1. Introduction

Consider a congestion system in which servers and customers may both be of more than one type. Such a system will be called a *heterogeneous congestion system*, and in such a system, the method used to assign customers to servers can have important consequences. For example, the method used to determine which heart attack victims are assigned to a coronary care unit can have a significant effect on the survival of heart attack victims (see Kao [5]).

In this paper we consider the problem of optimally assigning customers to servers in a two server heterogeneous congestion system with no waiting room. For other investigations into the problem of assigning customers to servers in heterogeneous congestion systems the reader is referred to Hall and Disney [3], Rolfe [8], Winston [11] and [12], Jarvis and Larson [4], and Carter, et al. [1].

### 2. Model Formulation

We assume that customers arrive according to a Poisson process with rate  $\lambda < \infty$ . The probability that an arrival is of type  $s$ ,  $-\infty < s < \infty$ , is  $P_s$ . We set  $\lambda_s = \lambda P_s$ , and define  $S = \{s : P_s > 0\}$ . An arrival who finds both servers idle may be assigned to either server; an arrival who finds only one server idle must be assigned to the idle server. Thus the only time a nontrivial decision must be made is when an arrival finds both servers idle. Upon assignment to server  $i$  ( $i = 1, 2$ ) a type  $s$  customer earns a reward  $R_{si}$ . Any type  $s$  customer who finds both servers busy is turned away and earns a reward  $\bar{R}_s$ . The distribution of a type  $s$  customer's service time is independent of the server to which he is assigned and has mean  $\mu_s$ . We assume that

$$\sup_{s \in S; i = 1, 2} |R_{si}| = M < \infty \quad \text{and} \quad (1)$$

\* All Notes are refereed.

† Accepted by Edward J. Ignall, former Departmental Editor; received September 1975. This paper has been with the author 3 months for 2 revisions.

‡ Indiana University.

$$\bar{\mu} = \sup_{s \in S} \mu_s < \infty \quad (2)$$

Carter, Chaiken and Ignall [1] consider a special case of our model in which service time is independent of customer type. Their Theorem 1 is a special case of our Theorem 2. Jarvis and Larson [4] consider the multiserver version of our model with  $S$  finite and service times exponential and independent of customer type. For an application of the above model to the utilization of coronary care units, see Kao [5]. Wrightson [14] has subsequently obtained results identical to our Theorem 3.

Our goal is to determine a method of assigning customers to servers that maximizes the average reward per unit time earned over an infinite horizon. To do so, we formulate the above problem as a semi-Markov decision process. We assume the reader is familiar with the standard terminology and results used to analyze such processes (see Ross [9]). For our problem, there is an obvious one to one correspondence between the set of stationary policies and the set  $\Delta$  of all subsets of  $S$ . For any  $T \in \Delta$ , we let  $T$  denote the stationary policy which assigns, whenever both servers are idle, customers of type  $s \in T$  to server 1 and all other customers to server 2.

A straightforward argument utilizing (1) and (2) shows that the mean time between successive epochs at which both servers are idle is bounded. Coupled with Theorem 6.19 of Ross [9], this yields the following result.

**THEOREM 1.** *There exists a stationary policy that maximizes the long-run average reward per unit time.*

### 3. Structure of the Optimal Policy

We now focus attention on the steady state behavior of the system when an arbitrary stationary policy,  $T$ , is followed. Given that  $T$  is followed, define

$P_{00}^T$  = the steady state probability that both servers are idle,

$P_{10}^T$  = the steady state probability that server 1 is occupied and server 2 is idle,

$P_{01}^T$  = the steady state probability that server 1 is occupied and server 2 is idle,

and  $P_{11}^T$  = the steady state probability that both servers are occupied. Since our choice of policy determines which, and not how many, servers are occupied,  $P_{00}^T = P_{00}$  and  $P_{11}^T = P_{11}$  for all  $T \in \Delta$ . It follows from [2] or [13] that  $P_{00} = 1/(1 + \rho + \rho^2/2)$  and  $P_{11} = P_{00}\rho^2/2$ , where  $\rho = \sum_S \lambda_s \mu_s$  and  $\sum_S$  is to be read as  $\sum_{s \in S}$ . To determine  $P_{ij}^T$ , we need to define, for arbitrary Borel sets  $\beta_1$  and  $\beta_2$ ,  $P^T(s_1, \beta_1, s_2, \beta_2, t)$  to be the probability that at time  $t$  server  $i$  has been occupied by a type  $s_i$  customer for the last  $t_i$  time units ( $t_i \in \beta_i$ ) given that at  $t = 0$  both servers were idle and policy  $T$  has been followed ( $s_i = 0$  denoting that server  $i$  has been idle for the last  $t_i$  time units). The argument used to prove Theorem 1 and the fact that arrivals are Poisson allows us to invoke Theorem 2.1 of Miller [7] to conclude that

$$\lim_{t \rightarrow \infty} P^T(s_1, \beta_1, s_2, \beta_2, t) = P^T(s_1, \beta_1, s_2, \beta_2), \quad (3)$$

then we have

$$P_{10}^T = \sum_S P^T(s, R^+, 0, R^+) \quad \text{and} \quad P_{01}^T = \sum_S P^T(0, R^+, s, R^+)$$

where  $R^+ = (0, \infty)$ . It is easy to see that the long-run average reward earned under policy  $T$  is independent of the initial state and is given by

$$\begin{aligned} \phi_T = & \sum_S P_{01}^T \lambda_s R_{s1} + \sum_S P_{10}^T \lambda_s R_{s2} \\ & + \sum_T \lambda_s P_{00}^T R_{s1} + \sum_{S-T} \lambda_s P_{00}^T R_{s2} + \sum_S \lambda_s \bar{R}_s P_{11}^T. \end{aligned} \quad (4)$$

To evaluate  $\phi_T$ , we need to evaluate  $P_{01}^T$  and  $P_{10}^T$ . Stidham's [9] version of  $L = \lambda W$  will be used. Focusing only on type  $s$  arrivals who are assigned (by choice or otherwise) to server 1, Theorem 1, equation 3, and Stidham's Corollary 10 allow us to conclude that

$$\bar{L}_s^T = \bar{\lambda}_s^T \bar{W}_s^T, \quad (5)$$

where  $\bar{L}_s^T = P^T(s, R^+, 0, R^+) + \sum_{s' \in S} P^T(s, R^+, s', R^+)$ ,  $\bar{\lambda}_s^T = \lambda_s [P_{01}^T + P_{00}^T]$ ,  $s \in T$ ,  $\bar{\lambda}_s^T = \lambda_s P_{01}^T$ ,  $s \in S - T$ , and  $\bar{W}_s^T = \mu_s$ .

Summing (5) over all  $s \in S$  yields

$$P_{10}^T + P_{11} = \rho P_{01}^T + \sum_T \lambda_s \mu_s P_{00}. \quad (6)$$

Since  $P_{00} + \dots + P_{11} = 1$ , (6) determines both  $P_{10}^T$  and  $P_{01}^T$ , a fact that will be used below in the proof of Theorem 3.

Define

$$\hat{T} = \left\{ s' : (R_{s'1} - R_{s'2}) / \mu_{s'} \geq \sum_S \lambda_s (R_{s1} - R_{s2}) / (1 + \rho) = K \right\}.$$

Then

**THEOREM 3.**  $\hat{T}$  is an optimal policy.

**PROOF.** By Theorem 1, it suffices to prove that  $\hat{T}$  is optimal among the class of stationary policies. Observe that for any  $T \in \Delta$  and  $s' \in S - T$  the fact the  $P_{00}$  and  $P_{11}$  are independent of  $T$  implies

$$P_{01}^{T \cup \{s'\}} - P_{01}^T = P_{10}^T - P_{10}^{T \cup \{s'\}}. \quad (7)$$

Call this quantity  $\Delta_{01}(T, s')$ . In combination with (7), subtraction of (6) for  $T \cup \{s'\}$  from (6) for  $T$  yields

$$\Delta_{01}(T, s') = -\lambda_{s'} \mu_{s'} P_{00} / (1 + \rho). \quad (8)$$

From (4),

$$\begin{aligned} \phi_{T \cup \{s'\}} - \phi_T = & \sum_S \lambda_s (R_{s1} - R_{s2}) \Delta_{01}(T, s) \\ & + \lambda_{s'} (R_{s'1} - R_{s'2}) P_{00}. \end{aligned} \quad (9)$$

Together (8) and (9) imply that  $\phi_{T \cup \{s'\}} > \phi_T$ , if and only if

$$\begin{aligned} P_{00} \sum_S \lambda_s (R_{s2} - R_{s1}) / 1 + \rho + P_{00} (R_{s'1} - R_{s'2}) / \mu_{s'} > 0 \\ \Leftrightarrow (R_{s'1} - R_{s'2}) / \mu_{s'} > \sum_S \lambda_s (R_{s1} - R_{s2}) / (1 + \rho) = K. \end{aligned} \quad (10)$$

We now let  $T$  denote an arbitrarily chosen optimal (stationary) policy, and proceed in a manner similar to the proof of Theorem 1 of [1]. If we let  $T_K = \{s : (R_{s1} - R_{s2}) / \mu_s = K\}$ , then we can complete the proof of the theorem by considering the following three cases:

- i.  $T \cap (S - \hat{T})$  is not empty.
- ii.  $(S - T) \cap (\hat{T} - T_K)$  is not empty.
- iii. Neither i nor ii is true and  $(S - T) \cap T_K$  is not empty.

If i is true, then for  $s \in T \cap (S - \hat{T})$ , (10) implies that  $\phi_{T-(s)} > \phi_T$ , thereby contradicting the optimality of  $T$ . If ii is applicable, then for  $s \in (S - T) \cap (\hat{T} - T_K)$ , (10) implies  $\phi_{T \cup (s)} > \phi_T$ , again contradicting the optimality of  $T$ . Therefore either iii is true or  $T = \hat{T}$ . If iii is true, then by (10),  $\phi_{T \cup T_K} > \phi_T$ , so  $T \cup T_K$  is optimal. Since  $T \cup T_K = \hat{T}$  whenever case iii is applicable, the proof of the theorem is complete.

#### 4. Remarks

Unfortunately, the above method of proof does not seem to be effective for systems having more than two servers. The reason for this is that for three or more servers, the analog of (7) does not supply sufficient information to enable us to successfully determine the structure of an optimal assignment policy.

If we assume that all  $R_{s1}$  are nonnegative and all  $R_{s2} = 0$  and identify assignment to server 1 with "acceptance" and assignment to server 2 with "rejection", then our model and results appear to be similar to the Streetwalker's Dilemma of Lippman and Ross [6] with Poisson arrivals; indeed the method of this paper enables one to successfully analyze the Streetwalker's Dilemma with Poisson arrivals. The main difference between the present model and the Streetwalker's Dilemma is that rejection of a customer in our model may force us to accept a future arrival, while in the Streetwalker's Dilemma rejection of an arrival places no constraints on the acceptance of future arrivals.<sup>1</sup>

<sup>1</sup> This work is an outgrowth of Chapter 7 of the author's Ph.D. thesis at Yale University. The author is grateful for the guidance of the members of his thesis committee, Matthew Sobel and Ward Whitt. This work was supported by the National Science Foundation and United States Public Health Service through grants GK-38121 and HS-00090-04, and the author is grateful for their financial support. The author also wishes to thank the departmental editor for several helpful suggestions.

#### References

1. CARTER, G., CHAIKEN, J. AND IGNALL, E., "Response Areas for Two Emergency Units," *Operations Research*, Vol. 20 (1972), pp. 571-594.
2. CHAIKEN, J. AND IGNALL, E., "An Extension of Erlang's Formulas which Distinguishes Individual Servers," *Journal of Applied Probability*, Vol. 9 (1972), pp. 192-197.
3. DISNEY, R. AND HALL, W., "Finite Queues in Parallel under a Generalized Channel Selection Rule," *Journal of Applied Probability*, Vol. 8 (1971), pp. 413-416.
4. JARVIS, J. AND LARSON, R., "Optimal Server Assignment Policies in M/M/N/O Queuing Systems with Distinguishable Servers and Customers Classes," Working Paper 06-74, Operations Research Center, Massachusetts Institute of Technology, 1974.
5. KAO, E., "Study of Patient Admission Policies for Specialized Care Facilities," Working Paper No. 29, Health Services Research Program, Yale University, 1973.
6. LIPPMAN, S., AND ROSS, S., "The Streetwalker's Dilemma: A Job Shop Model," *SIAM J. Appl. Math.*, Vol. 20 (1971), pp. 336-344.
7. MILLER, E., "Existence of Limits in Regenerative Processes," *Ann. Math. Statist.*, Vol. 43 (1972), pp. 1275-1282.
8. ROLFE, A., "The Control of Multiple Facility, Multiple Channel Queuing System with Parallel Input Streams," Technical Report No. 22, Graduate School of Business, Stanford University, 1965.
9. ROSS, S., *Applied Probability with Optimization Applications*, Holden-Day, 1970.
10. STIDHAM, S. Jr., "L =  $\lambda W$ : A Discounted Analogue and a New Proof," *Opns. Res.*, Vol. 20 (1972), pp. 1115-1125.
11. WINSTON, W., "Assignment of Customers to Servers in a Heterogeneous Queuing with Switching," *Operations Research*, Vol. 25 (1977), pp. 469-483.
12. ———, "Optimal Dynamic Rules for Assigning Customers to Servers in a Heterogeneous Queuing System," *Naval Research Logistics Quarterly*, Vol. 24 (June 1977), pp. 273-300.
13. WOLFF, R. AND WRIGHTSON, C., "An Extension of Erlang's Loss Formula," *Journal of Applied Probability*, Vol. 13 (1976), pp. 628-632.
14. WRIGHTSON, C., "The Design of Response Areas for Emergency Service Units," paper presented at ORSA-TIMS meeting, May 1977.

Copyright 1978, by INFORMS, all rights reserved. Copyright of Management Science is the property of INFORMS: Institute for Operations Research and its content may not be copied or emailed to multiple sites or posted to a listserv without the copyright holder's express written permission. However, users may print, download, or email articles for individual use.