# Approximate Analysis of Optimal Operating Policies for a GI/G/1 Queueing System

By

Toshikazu Kimura*, Katsuhisa Ohno* and Hisashi Mine*

### Abstract

In this paper, applying a method of diffusion approximation, we consider optimal operating policies for a $GI/G/1$ queueing system with a removable server. The following costs are incurred in the system: a cost per unit time of keeping the server running, fixed costs for turning the server on or off, and a holding cost per customer in the system per unit time. The average cost rate is used as a criterion for optimality. By using a couple of diffusion processes that approximate the number of customers in the system, an explicit form of the average cost rate is derived. Furthermore, some sufficient conditions under which the optimal operating policy falls into specific forms are obtained. It is examined numerically how the boundary condition at the origin of the diffusion process effects the optimal operating policy and its cost.

## 1. Introduction

Queueing systems with a removable server have been studied by many researchers, e.g., Yadin and Naor[16], Heyman[8], Sobel[14] and Bell[2]. In this paper, such a system under the general conditions is analyzed by using diffusion approximation, and some useful properties of optimal policies for operating the system are provided. In the standard $GI/G/1$ queueing system, the server always serves arriving customers, i.e., it is always *on* so far as the system is not empty. However, the server in the present $GI/G/1$ system is allowed to remain dormant in front of arriving customers, i.e., it may remain *off*. Furthermore, the following various costs are incurred in the system. First, costs for keeping the server off or on are incurred per unit time. It is assumed that the cost for keeping the server on (*running cost*) is greater than that for keeping it off (*dormant cost*). This is an incentive to turn the server off whenever the system becomes empty. Secondly, fixed costs for turning the server on and turning it off are

---

* Department of Applied Mathematics and Physics

incurred. The former is called a *start-up cost* and the latter a *shut-down cost*. The existence of the start-up cost is an incentive to keep the server off even if customers are waiting in the queue. In addition to the above costs there is a cost per unit time as a penalty for the delay of customers. This cost is proportional to the number of customers in the system, and is called a *holding cost*.

The problem of controlling the system concerns what policy would be optimal for turning the server off or on with respect to the information about the number of customers, and whether or not the server is off. As a criterion for optimality, we use the average cost per unit time over an infinite horizon. This problem for the $M/G/1$ system has been first studied by Heyman[8]. For the $GI/G/1$ system with a more general cost structure, Sobel[14] obtained a sufficient condition under which the optimal stationary policy is, so called, an $(n, N)$ policy $(0 \leq n \leq N)$. He mentioned that $n=0$ in the queueing system with the above cost structure. Moreover, Yadin and Naor[16] investigated the system in which both the running cost and the dormant cost vanish, and the operations for the start-up and the shut-down require stochastically distributed times. Their operating policy is also the $(0, N)$ policy.

In connection with control problems of dams and storage systems, much extensive research on a controlled Brownian motion process has been recently carried out[1,3~5,7,12]. Rath[12] investigated a controlled diffusion process with the same cost structure and optimality criterion as in this paper. Approximating the underlying diffusion process by a sequence of random walk and applying the theory of Markov decision processes, he showed that the optimal policy is the $(n, N)$ policy. He obtained, however, few results about the properties of optimal policies, since the diffusion parameters are not specified in terms of the parameters of the original controlled queueing problem. Chernoff and Petkau[3] proposed a slight generalization of Rath's results in which a non-linear holding cost is allowed. Their computational procedure for the optimal policy is, however, so complicated that one must solve a system of non-linear equations even for the linear cost structure. A Brownian motion control problem with quadratic cost was analyzed by Doshi[4].

In this paper, considering these results, we shall investigate both the properties of the optimal operating policies and the accuracy of the approximation for the controlled queueing system with a removable server. Several assumptions and notation specifying the system are introduced in Section 2. In Section 3, using diffusion processes approximating the number of customers in the system, we analyze the behaviour of the system under the $(0, N)$ policy, and obtain an explicit expression for the average cost rate. Section 4 deals with the characterization of the optimal operating policy. It is shown that the theorems in Section 4 are rather new results which also hold exactly for the $M/G/1$ system. Section 5 examines numerically how the boundary condition at the origin of the diffusion process effects the optimal operating policy and its cost.

## 2. Assumptions and Notation

Let us define *idle periods* as time intervals during which the server is off. Time intervals when the server is on are called *busy periods*. Moreover, a successive idle and busy period is called a *busy cycle*. It is noted that the definitions of idle and busy periods are different from the usual queueing terminology in which these periods are determined by whether or not the system is empty.

Assumptions on arrival and service processes are:

(a) *Arrival process*

Let $u_n$ ($n=1, 2, \cdots$) denote the interarrival time between the $(n-1)^{st}$ and $n^{th}$ arriving customers; $\{u_n, n=1, 2, \cdots\}$ is a sequence of nonnegative i.i.d. random variables with $E[u_n]=1/\lambda$ ($<\infty$) and $\mathrm{Var}[u_n]=\sigma_a^2$ ($<\infty$).

(b) *Service process*

Let $v_n$ ($n=1, 2, \cdots$) denote the service time of the $n^{th}$ arriving customer; $\{v_n, n=1, 2, \cdots\}$ is a sequence of nonnegative i.i.d. random variables with $E[v_n]=1/\mu$ ($<\infty$) and $\mathrm{Var}[v_n]=\sigma_s^2$ ($<\infty$). It is assumed that the traffic intensity $\rho\equiv\lambda/\mu$ is less than unity. Moreover, the sequence $\{v_n, n=1, 2, \cdots\}$ is independent of $\{u_n, n=1, 2, \cdots\}$. The order of servicing customers is arbitrary as far as it is work-conserving.

(c) *Cost structure*

Notations used for the costs incurred in the system are as follows:

$r_1$: dormant cost rate [$/hour],

$r_2$: running cost rate [$/hour],

$R_1$: start-up cost [$],

$R_2$: shut-down cost [$],

$h$: holding cost rate [$/customer·hour].

In this linear cost structure, all costs are assumed to be nonnegative and finite. Further, it is assumed that $h>0$ and $r_2\geq r_1$.

(d) *Criterion for optimality*

We adopt an average cost per unit time over an infinite horizon as the cost function to be minimized. Let $C(t)$ be the total cost incurred during the time interval $(0, t]$. Then the optimal policy is a policy that minimizes

$$\lim_{t\to\infty} \frac{1}{t} E[C(t)]. \tag{1}$$

(e) *Operating policies*

The forms of operating policies are closely related to the choice of decision epochs.

For the $M/G/1$ system, Heyman[8] proved that the optimal policy has one of the following forms if the server should be turned on at any arrival epoch, or off at any service completion epoch:

(i)  Always turn the server on, or never turn the server off.

(ii)  Turn the server on when $N$ ($\geq 1$) or more customers are present, turn the server off when the system is empty. It is called a $(0, N)$ policy.

(iii)  Always turn the server off, or never turn the server on. It may be called a $(0, \infty)$ policy.

For the $GI/G/1$ system, Sobel[14] showed, under weaker conditions for the costs and the same decision epochs, that the optimal policy is an $(n, N)$ policy. He pointed out that $n=0$ under the linear cost structure. Consequently, following their results, we can restrict the form of operating policies to

$\pi_0$:  a policy that the server is always turned on,

and

$\pi_N$:  a $(0, N)$ policy ($1 \leq N \leq \infty$).

In the formulation of Section 3, the decision epochs do not directly correspond to the arrival or service completion epochs since diffusion processes have continuous state spaces. Rath[12] proved, however, that the optimal policy for controlled diffusion processes with the same cost structure is also the $(n, N)$ policy. For the assumed queueing contents, the value of $n$ must be zero. Thus, the choice of the decision epochs does not have any effect on the subsequent analysis.

## 3.  Formulation by a Diffusion Model

In this section we shall explicitly derive the average cost rate under the policy $\pi_N$ by using two appropriate diffusion processes. Since the average cost rate does not depend on the initial state, it is assumed, without loss of generality, that an idle period begins at time zero. Let $Q(t)$ denote the number of customers in the system at time $t$, and consider the behaviour of the path $Q(t)$ in one busy cycle under the policy $\pi_N$. The sample path $Q(t)$ is approximated by two independent diffusion processes $\{Y(t);$ $t \geq 0\}$ and $\{Z(t); t \geq 0\}$. The former approximates $Q(t)$ in the idle period and the latter in the busy period. In the idle period, $Q(t)$ increases monotonically according to the arrival of customers because the server is never turned on during that period. Let $A(t)$ denote the total number of arrivals in the time interval $(0, t]$. Since $Q(0)=0$, $Q(t)=A(t)$ in the idle period. Hence, using the results in [9], the path $Q(t)$ in the idle period can be approximated by a diffusion process $\{Y(t); t \geq 0\}$ with the infinitesimal mean $b_1$ and the infinitesimal variance $a_1$, where

$$b_1 \equiv \lim_{\Delta t \downarrow 0} \frac{E[Y(t+\Delta t)-Y(t)\,|\,Y(t)]}{\Delta t} = \lambda, \tag{2}$$

and

$$a_1 \equiv \lim_{\Delta t \downarrow 0} \frac{\mathrm{Var}\,[Y(t+\Delta t)-Y(t)\,|\,Y(t)]}{\Delta t} = \lambda^3 \sigma_a{}^2. \tag{3}$$

When the path first reaches level $N$, the busy period starts and the server is turned on. Let $D(t)$ denote the total number of departures in the time interval $(0,\ t]$. Then it follows that $Q(t)=A(t)-D(t)$ in the busy period. That is, the path $Q(t)$ in the busy period behaves similarly to that of the usual $GI/G/1$ system. Hence, it can be approximated by the diffusion process $\{Z(t);\ t \geq 0\}$ with

$$b_2 \equiv \lim_{\Delta t \downarrow 0} \frac{E[Z(t+\Delta t)-Z(t)\,|\,Z(t)]}{\Delta t} = \lambda - \mu, \tag{4}$$

and

$$a_2 \equiv \lim_{\Delta t \downarrow 0} \frac{\mathrm{Var}\,[Z(t+\Delta t)-Z(t)\,|\,Z(t)]}{\Delta t} = \lambda^3 \sigma_a{}^2 + \mu^3 \sigma_s{}^2. \tag{5}$$

From the physical meaning of $Q(t)$, the processes $Y(t)$ and $Z(t)$ should have the regions $[0,\ N]$ and $[0,\ \infty)$, respectively. However, it is assumed for analytical convenience that the former has the region $(-\infty,\ N]$.

For these diffusion processes, define the stopping times as follows:

$$T_1(x_0, N) \equiv \inf\,\{t \geq 0\,|\,Y(t)=N,\ Y(0)=x_0\}, \tag{6}$$

and

$$T_2(x_0, 0) \equiv \inf\,\{t \geq 0\,|\,Z(t)=0,\ Z(0)=x_0\}. \tag{7}$$

The stopping times $T_1 \equiv T_1(0,\ N)$ and $T_2 \equiv T_2\,(N,\ 0)$ approximate the lengths of the idle and busy periods, respectively. From $b_1 > 0$, we have $E[T_1] < \infty$. The assumption $\rho < 1$ implies that $b_2 < 0$ and hence $E[T_2] < \infty$. Thus, in an average sense, one busy cycle terminates at the finite time $t = T_1 + T_2$.

The following lemma is directly derived from Theorem 7.5 in [13, p. 160].

**Lemma 1.** *Let $V(N)$ denote the average cost rate of the system under the policy* $\pi_N$, i.e.,

$$V(N) = \lim_{t \to \infty} \frac{1}{t} E[C(t)]. \tag{8}$$

*Then, for* $N = 1,\ 2,\ \cdots,$

$$V(N) = \frac{1}{E[T_1]+E[T_2]} \Big\{ E\Big[ \int_0^{T_1} c_1(Y(u))du \Big] \\ + E\Big[ \int_0^{T_2} c_2(Z(u))du \Big] + R_1 + R_2 \Big\}, \tag{9}$$

*where*

$$c_i(x) = hx + r_i, \qquad i = 1, 2. \tag{10}$$

In order to evaluate the average cost rate given by (9), it remains to calculate both the mean lengths of the idle and busy periods and the expected costs incurred during these periods. The next lemma provides the differential equations that the expected costs satisfy.

**Lemma 2.** *Let*

$$v_1(x_0) = E\left[ \int_0^{T_1(x_0, N)} c_1(Y(u)) du \right]$$

*and*

$$v_2(x_0) = E\left[ \int_0^{T_2(x_0, 0)} c_2(Z(u)) du \right].$$

*Then, $v_i(x_0)$ satisfies the ordinary differential equation*

$$\frac{1}{2} a_i \frac{d^2 v_i}{dx_0^2} + b_i \frac{dv_i}{dx_0} + c_i(x_0) = 0, \qquad 0 < x_0 < N, \tag{11}$$

*for $i = 1, 2$. Their boundary conditions are given by $v_1(N) = 0$ and $v_2(0) = 0$, respectively.*

**Proof:** Using Theorem 1 in [11, p. 149], we can easily derive (11). The boundary conditions correspond to imposing absorbing barriers at $x = N$ for the idle period and $x = 0$ for the busy one. ☐

The mean lengths of the idle and busy periods can be obtained from Lemma 2, as shown in the next lemma.

**Lemma 3.** *Let $m_1(x_0) = E[T_1(x_0, N)]$ and $m_2(x_0) = E[T_2(x_0, 0)]$. Then, $m_i(x_0)$ satisfies the ordinary differential equation*

$$\frac{1}{2} a_i \frac{d^2 m_i}{dx_0^2} + b_i \frac{dm_i}{dx_0} = -1, \qquad 0 < x_0 < N, \tag{12}$$

*for $i = 1, 2$. Their boundary conditions are given by $m_1(N) = 0$ and $m_2(0) = 0$, respectively.*

**Proof:** Substituting $c_i(x_0) \equiv 1$ in Lemma 2 immediately leads to the result. ☐

Thus, the average cost rate $\{V(N)\}$ is obtained from these lemmas.

**Theorem 1.** *Let $V(N)$ be the average cost rate of the system under the policy $\pi_N$. Then*

$$V(0) = r_2 + hL(1), \tag{13}$$

*and for $N = 1, 2, \cdots$,*

$$V(N) = r_1 + \rho(r_2 - r_1) + hL(N) + \frac{1}{N} \lambda(1 - \rho)(R_1 + R_2), \tag{14}$$

*where*

$$L(N) = \frac{1}{2}\Big[N + \frac{\rho}{1-\rho}\{\mu^2\sigma_s{}^2 + \lambda(2\lambda-\mu)\sigma_a{}^2\}\Big].$$  (15)

**Proof:**  The proof is directly done by Lemmas 1∼3.  Solving (11) with the given boundary conditions and setting $x_0$ to each initial value, we have

$$v_1(0) = \frac{h}{2b_1}N^2 - \frac{1}{b_1}\Big(\frac{a_1 h}{2b_1} - r_1\Big)N,$$  (16)

and

$$v_2(N) = -\frac{h}{2b_2}N^2 + \frac{1}{b_2}\Big(\frac{a_2 h}{2b_2} - r_2\Big)N.$$  (17)

The mean lengths of the periods are derived from Lemma 3, or by substituting $r_1 = r_2 \equiv 1$ and $h \equiv 0$ into (16) and (17).  That is,

$$m_1(0) = \frac{N}{b_1},$$  (18)

and

$$m_2(N) = -\frac{N}{b_2}.$$  (19)

Then, by applying Lemma 1 and using (2)∼(5) and (16)∼(19), it will yield (14).  From the structure of (14), it is considered that $L(N)$ represents the mean number of customers in the steady state under the policy $\pi_N$.  That is, $L(1)$ must agree with the mean number of customers in the usual sense.  Hence, the average cost rate using the policy $\pi_0$ is clearly given by (13).  □

The average cost rate (14) has a very similar form as the exact one for the $M/G/1$ system[8].  That is, the terms, except for $hL(N)$, completely agree with those of the exact results.  For the $M/G/1$ system, the term $L(N)$ can be rewritten as

$$L(N) = \frac{1}{2}(N-1) + \frac{\rho(\mu^2\sigma_s{}^2+1)}{2(1-\rho)}.$$  (20)

Let $c_s$ be the coefficient of variation of service times.  Then, for $c_s{}^2 = 1$, e.g., for the $M/M/1$ system, $L(N)$ becomes equivalent to the exact one.  For other values of $c_s{}^2$, it follows from [16] that $L(N)$ has an error evaluated by

$$\frac{\rho(\mu^2\sigma_s{}^2+1)}{2(1-\rho)} - \Big\{\frac{\rho^2+\lambda^2\sigma_s{}^2}{2(1-\rho)} + \rho\Big\} = \frac{\rho(c_s{}^2-1)}{2}.$$  (21)

Hence, $V(N)$ overestimates or underestimates the exact average cost rate, according as $c_s{}^2 > 1$ or $c_s{}^2 < 1$.  The relative error of (21), however, decreases as the traffic becomes heavy, since the error (21) increases with the linear order of $\rho$, whereas the number of customers increases with the order of $(1-\rho)^{-1}$.

## 4. Optimal Operating Policies

The optimal operating policy that minimizes the average cost rate derived in Section 3 will be found in this section. Although the average cost rates, $\{V(N), N \geqq 1\}$, are defined only on the set of the natural number, we extensively regard $V(N)$ as a function of the positive real number $N$ for the sake of analytical convenience. Then, by differentiating $V(N)$ by $N$, it is shown that $V(N)$ has a unique minimum at $N = \bar{N}$ because of its convexity, where

$$\bar{N} = \sqrt{\frac{2\lambda(1-\rho)(R_1+R_2)}{h}} . \tag{22}$$

Hence, the best positive integer value of $N$ is given by one of two integer points adjacent to $\bar{N}$. Consequently, the optimal value of $N$ can be determined by comparing the values of $V(N)$ at these integers and $V(0)$. For the $M/G/1$ system, Heyman[8] obtained the same result as above. Further, it also agrees with the result of Yadin and Naor[16] in a case where the switching times of start-up and shut-down are ignored.

The next two theorems and corollaries provide some sufficient conditions which put the optimal operating policy in specific forms. It should be noted that these new theorems hold exactly for the $M/G/1$ system, because the error of $V(N)$ is independent of $N$.

**Theorem 2.** *If $r_1 = r_2$, then the optimal operating policy is $\pi_0$.*

**Proof:** The proof is executed by distinguishing two cases;

*Case 1.* $\bar{N} \leqq \frac{1}{2}$. For this case it is clear that $\min_{N} V(N) = V(1)$, where the minimum operation is chosen over the set of the natural number. Hence,

$$\min_{N} V(N) - V(0) = V(1) - V(0)$$
$$= \lambda(1-\rho)(R_1+R_2) > 0.$$

That is, the optimal operating policy is $\pi_0$.

*Case 2.* $\bar{N} > \frac{1}{2}$. Since $V(\bar{N})$ must be less than or equal to $\min_{N} V(N)$, it follows from (13)~(15) and (22) that

$$\min_{N} V(N) - V(0) \geqq V(\bar{N}) - V(0)$$
$$= h\{L(\bar{N}) - L(1)\} + \frac{1}{2} h\bar{N}$$
$$= h\left(\bar{N} - \frac{1}{2}\right) > 0.$$

That is, the optimal operating policy is $\pi_0$. □

From the result of Theorem 2, we assume hereafter that $r_2 > r_1$.

## Theorem 3.

(i)  *If* $r_2-r_1<\mu(R_1+R_2)$, *then there exists a unique* $\lambda^*\in(0,\mu)$, *such that for any* $\lambda\in$ $[\lambda^*,\mu)$ *the optimal operating policy is* $\pi_0$.

(ii)  *Otheriwse, there exists a unique* $\lambda^{**}\in(0,\mu)$, *such that for any* $\lambda\in[\lambda^{**},\mu)$ *the optimal operating policy is* $\pi_1$.

**Proof:**  From (22), we have

$$\bar{N}^2=\frac{2(R_1+R_2)}{\mu h}\lambda(\mu-\lambda). \tag{23}$$

That is, $\bar{N}^2$ can be regarded as a quadratic function of $\lambda$.  Therefore, $\bar{N}^2$ is monotonically decreasing for $\lambda>\mu/2$ and gets less than unity as $\lambda$ increases, i.e., $2(R_1+R_2)$ $\cdot\lambda(\mu-\lambda)/\mu h\leq1$.  This inequality can be rewritten as

$$\left(\lambda-\frac{\mu}{2}\right)^2\geq\frac{\mu^2}{4}\left\{1-\frac{2h}{\mu(R_1+R_2)}\right\}. \tag{24}$$

Hence, if $\mu(R_1+R_2)\leq2h$, then (24) always holds, that is, $\bar{N}\leq1$ for any $\lambda\in[0,\mu)$.  On the other hand, if $\mu(R_1+R_2)>2h$, then $\bar{N}\leq1$ for $\lambda\geq\lambda_1$ with

$$\lambda_1=\frac{\mu}{2}\left\{1+\sqrt{1-\frac{2h}{\mu(R_1+R_2)}}\right\}. \tag{25}$$

It is clear that $\lambda_1\in(\mu/2,\mu)$.  Consequently, a sufficient condition for $\bar{N}\leq1$, which is independent of the costs, is $\lambda\geq\lambda_1$.  Assume here that $\lambda$ is sufficiently large so that $\bar{N}\leq1$.  Then,

$$\min_N V(N)-V(0)=V(1)-V(0)$$
$$=(1-\rho)\{\lambda(R_1+R_2)-(r_2-r_1)\}.$$

Hence, if $\lambda\geq\lambda_2\equiv(r_2-r_1)/(R_1+R_2)$ and $0<\lambda_2<\mu$, it follows that $\min_N V(N)\geq V(0)$ and the optimal operating policy is $\pi_0$.  In other words, if $r_2-r_1<\mu(R_1+R_2)$, then the optimal operating policy is $\pi_0$ for $\lambda\geq\max(\lambda_1,\lambda_2)\in(0,\mu)$.  For a case where $\lambda_2\geq\mu$, i.e., $r_2-r_1\geq\mu(R_1+R_2)$, it follows for any $\lambda$ with $\bar{N}\leq1$ that

$$\min_N V(N)=V(1)<V(0),$$

whereby the optimal operating policy is $\pi_1$ for $\lambda$ satisfying $\bar{N}\leq1$.  Such $\lambda$ is given, for example, by $\lambda\in[\lambda_1,\mu)$.  Thus the proof is completed.  $\square$

**Corollary 1.**  *If* $r_2-r_1<\mu(R_1+R_2)\leq2h$, *then the optimal operating policy is*

$$\pi_1 \quad for \quad 0<\lambda<(r_2-r_1)/(R_1+R_2),$$

*and*

$$\pi_0 \quad for \quad (r_2-r_1)/(R_1+R_2)\leq\lambda<\mu.$$

**Corollary 2.** *If* $r_2-r_1\geq\mu(R_1+R_2)$ *and* $2h\geq\mu(R_1+R_2)$, *then the optimal operating policy is* $\pi_1$.

The proofs of these corollaries obviously follow that of Theorem 3, and hence are omitted here.

Theorem 3 and its corollaries mean that the optimal operating policy eventually falls into $\pi_0$ or $\pi_1$ as the traffic becomes heavy.

**Remark 1.** Although it has been assumed that $h\neq 0$, a case where $h=0$ could occur in a system where the holding cost is ignored or included in the running cost. For this case, it is clear that the optimal operating policy is $\pi_\infty$.

## 5. Boundary Conditions at the Origin

Since $Q(t)$ represents the number of customers in the system, it cannot take a negative value. However, no restriction is imposed on the process $Y(t)$ that approximates $Q(t)$ in the idle period. With respect to the average behaviour of $Y(t)$, it follows from $b_1>0$ that

$$E[Y(t+\Delta t)-Y(t)]>0, \qquad \text{for } \Delta t>0. \tag{26}$$

Hence, it is guaranteed in the average sense that the process $Y(t)$ increases monotonically. Each trial of the path $Y(t)$, however, has the possibility of violating the nonnegative restriction. Specifically, for a case where the variance of interarrival times $\sigma_a^2$ is large, the probability that the path $Y(t)$ invades the negative region becomes large since the diffusion coefficient $a_1$ is proportional to $\sigma_a^2$. Considering this possibility, Whitt[15] adopted the process $Y(t)$ with a reflecting barrier at the origin. His results can be obtained easily by modifying the boundary conditions in Lemmas 2 and 3 as follows [10]:

$$\frac{dv_1}{dx_0}\bigg|_{x_0=0}=0, \tag{27}$$

and

$$\frac{dm_1}{dx_0}\bigg|_{x_0=0}=0. \tag{28}$$

Solving the differential equations (11) and (12) with the boundary conditions (27) and (28), respectively, and setting $x_0=0$, we have

$$\hat{v}_1(0)=\frac{h}{2b_1}N^2-\left(\frac{a_1h}{2b_1}-r_1\right)E[\hat{T}_1], \tag{29}$$

and

$$E[\hat{T}_1]=\frac{N}{b_1}+\frac{a_1}{2b_1^2}\{\exp(-2b_1N/a_1)-1\}, \tag{30}$$

where $\hat{\ }$ denotes the characteristics of the process $Y(t)$ with a reflecting barrier at the origin. Hence, the average cost rate under the policy $\pi_N(N \geqq 1)$ is given by

$$\hat{V}(N) = \frac{1}{E[\hat{T}_1] + E[T_2]} \{\hat{v}_1(0) + v_2(N) + R_1 + R_2\}$$

$$= r_1 + \hat{\rho}(r_2 - r_1) + h\hat{L}(N) + \frac{R_1 + R_2}{E[\hat{T}_1] + E[T_2]}, \tag{31}$$

where

$$\hat{\rho} = \frac{E[T_2]}{E[\hat{T}_1] + E[T_2]}, \tag{32}$$

and

$$\hat{L}(N) = \frac{1}{2(E[\hat{T}_1] + E[T_2])} \left\{ \left( \frac{1}{b_1} - \frac{1}{b_2} \right) N^2 - \frac{a_1}{b_1} E[\hat{T}_1] - \frac{a_2}{b_2} E[T_2] \right\}. \tag{33}$$

Furthermore, from the statements in Section 3, the average cost rate under the policy $\pi_0$ is given by

$$\hat{V}(0) = r_2 + h\hat{L}(1). \tag{34}$$

**Remark 2.** The right hand side of (31) corresponds to that of (14). The symbol $\hat{\rho}$ represents the utilization of the system, and the term $(E[\hat{T}_1] + E[T_2])^{-1}$ represents the number of busy cycles per unit time. It follows from (19), (30) and (32) that

$$\lim_{N \to \infty} \hat{\rho} = \rho. \tag{35}$$

**Remark 3.** If the variance of interarrival times $\sigma_a{}^2$ vanishes, that is, if the system is $D/G/1$, then $\{\hat{V}(N)\}$ agrees with $\{V(N)\}$. Therefore, the results in Sections 3 and 4 also hold for the approximate solutions with the reflecting barrier. This fact is partially noted in [15].

It follows from the above remarks that a boundary condition at the origin for the process $Y(t)$ effects the solutions to some degree. Hereafter, taking account of the importance of the boundary condition, we deal with a new and more natural boundary condition, which is relevant to the average behaviour of $Y(t)$. Let $T_0$ be the time from the beginning of one busy cycle to the first arrival of a customer in this busy cycle. Then it is known that the relation $E[T_0] = 1/\lambda$ does not always hold. Moreover, it is difficult in general to give an exact expression of $E[T_0]$. We assume in the following that $T_0$ is a stationary residual life time of the interarrival times $\{u_n\}$. Thus, its mean is given by

$$E[T_0] = \frac{\lambda^2 \sigma_a{}^2 + 1}{2\lambda}. \tag{36}$$

This assumption seems to be plausible because what effects the average cost rate is only relevant to the stationary behaviour of the system. Since the first passage time $T_1(0,$

$N$) is a sum of $T_0$ and $T_1$ (1, $N$), it is quite natural to adopt the following boundary condition at the origin:

$$m_1(0) = E[T_0] + m_1(1). \tag{37}$$

As for the costs incurred during the idle period, the boundary condition is similarly given by

$$v_1(0) = r_1 E[T_0] + v_1(1), \tag{38}$$

because the cost incurred in the time interval $T_0$ is the dormant cost only. Taking the boundary condition (38) into account in Lemma 2, we obtain

$$\tilde{v}_1(0) = \frac{h}{2b_1} N^2 - \frac{1}{b_1} \left( \frac{a_1 h}{2b_1} - r_1 \right) N$$
$$+ \frac{1 - \exp(-2b_1 N / a_1)}{1 - \exp(-2b_1 / a_1)} \left\{ r_1 \left( E[T_0] - \frac{1}{b_1} \right) + \frac{h}{2b_1} \left( \frac{a_1}{b_1} - 1 \right) \right\}. \tag{39}$$

Substituting $r_1 = 1$ and $h = 0$ in (39) leads to

$$\tilde{m}_1(0) = \frac{N}{b_1} + \frac{1 - \exp(-2b_1 N / a_1)}{1 - \exp(-2b_1 / a_1)} \left( E[T_0] - \frac{1}{b_1} \right), \tag{40}$$

where $\tilde{\ }$ denotes the characteristics of the process $Y(t)$ with the boundary condition (37) or (38). These boundary conditions physically correspond to the process which remains at the origin for the time interval $T_0$, and then jumps instantaneously to $x = 1$, and thereafter starts from scratch. If the arrival process is a Poisson process, the process $Y(t)$ is called an *elementary return process*, and has been investigated by Feller[6]. Moreover, since $a_1 = b_1 = E[T_0]^{-1} = \lambda$ for the Poisson arrival case, we have $\tilde{v}_1(0) = v_1(0)$ and $\tilde{m}_1(0) = m_1(0)$.

In order to examine the accuracy of the diffusion approximation and to show the

Table 1.  *The Optimal Operating Policy $\pi_N*$ and Its Cost Rate for the $M/E_2/1$ System*
($r_1 = 10$, $r_2 = 50$, $R_1 = R_2 = 100$, $h = 1$, $\mu^{-1} = 1.0$).

| $\rho$ | exact | | diffusion approximation | | | | | |
|---|---|---|---|---|---|---|---|---|
| | $N*$ | cost rate | $N*$ | $V(N*)$ | relative error (%) | $N*$ | $\hat{V}(N*)$ | relative error (%) |
| 0.1 | 6 | 19.608 | 6 | 19.583 | (−0.127) | 7 | 20.352 | (3.794) |
| 0.2 | 8 | 25.737 | 8 | 25.678 | (−0.194) | 9 | 26.498 | (2.957) |
| 0.3 | 9 | 31.063 | 9 | 30.988 | (−0.241) | 10 | 31.801 | (2.378) |
| 0.4 | 10 | 35.900 | 10 | 35.800 | (−0.278) | 11 | 36.602 | (1.957) |
| 0.5 | 10 | 40.375 | 10 | 40.250 | (−0.309) | 11 | 41.011 | (1.576) |
| 0.6 | 10 | 44.575 | 10 | 44.425 | (−0.366) | 11 | 45.136 | (1.259) |
| 0.7 | 9 | 48.591 | 9 | 48.416 | (−0.360) | 10 | 49.043 | (0.929) |
| 0.8 | 8 | 52.700 | 8 | 52.500 | (−0.379) | 9 | 53.039 | (0.643) |
| 0.9 | 0 | 56.975 | 0 | 56.750 | (−0.394) | 0 | 57.077 | (0.180) |

Table 2.  *The Optimal Operating Policy $\pi_N{}^*$ and Its Cost Rate for the $E_2/E_5/1$ System*
($r_1=10$, $r_2=50$, $R_1=R_2=100$, $h=1$, $\mu^{-1}=1.0$).

| $\rho$ | simulation | | diffusion approximation | | | | | |
|---|---|---|---|---|---|---|---|---|
| | $N^*$ | cost rate | $N^*$ | $V(N^*)$ | $N^*$ | $\hat{V}(N^*)$ | $N^*$ | $\tilde{V}(N^*)$ |
| 0.1 | 6 | 19.998 | 6 | 19.788 | 6 | 20.180 | 6 | 20.137 |
| 0.2 | 8 | 26.005 | 8 | 25.837 | 8 | 26.250 | 8 | 26.225 |
| 0.3 | 10 | 31.175 | 9 | 31.066 | 10 | 31.480 | 10 | 31.464 |
| 0.4 | 10 | 35.895 | 10 | 35.783 | 10 | 36.179 | 10 | 36.167 |
| 0.5 | 10 | 40.479 | 10 | 40.100 | 10 | 40.484 | 10 | 40.475 |
| 0.6 | 10 | 44.250 | 10 | 44.075 | 10 | 44.421 | 10 | 44.415 |
| 0.7 | 9 | 47.924 | 9 | 47.733 | 9 | 48.052 | 9 | 48.047 |
| 0.8 | 8 | 51.500 | 8 | 51.150 | 8 | 51.410 | 8 | 51.407 |
| 0.9 | 0 | 53.489 | 0 | 53.400 | 0 | 53.491 | 0 | 53.461 |

effects of its boundary conditions, the approximate results in Sections 3 and 5 are compared numerically with the known exact results. Table 1 shows the exact and approximate solutions of the optimal operating policies and their costs for the $M/E_2/1$ system. Note that the cost structure satisfies the condition (i) of Theorem 3. No significant errors are found either for the optimal operating policy or for the corresponding cost rate. It is shown from many other numerical results for the $M/G/1$ system that the approximation $\{V(N)\}$ ($=\{\tilde{V}(N)\}$) is more accurate than $\{\hat{V}(N)\}$, except for the heavy traffic.

For general arrival processes, it is difficult to obtain the exact results analytically. Hence, we shall compare the approximate results with the results obtained from the GPSS simulation. These results for the $E_2/E_5/1$ system are shown in Table 2. It follows from this table that the approximation $\{\tilde{V}(N)\}$ behaves more like $\{\hat{V}(N)\}$ than $\{V(N)\}$, and that $\{\tilde{V}(N)\}$ and $\{\hat{V}(N)\}$ are slightly more accurate than $\{V(N)\}$. This tendency becomes apparent as the coefficient of variation of interarrival times becomes small. Consequently, it seems that the approximation $\{\tilde{V}(N)\}$ is best suited for all cases of the $GI/G/1$ system.

## 6.  Concluding Remarks

This paper characterizes the optimal operating policy that minimizes the average cost rate. As another criterion for optimality, the expected total discounted cost, $E\left[\int_0^\infty e^{-\beta t}dC(t)\right]$, is also frequently used for evaluating the costs incurred in the system[2], where $\beta(>0)$ denotes a discount rate. In particular, this discounted cost seems to be appropriate for investment problems. By using the same diffusion processes as defined in Section 3, it is possible to calculate the discounted cost and to analyze an optimal operating policy.

## References

1) J. A. Bather, "A Diffusion Model for the Control of a Dam," *Journal of Applied Probability*, Vol. 5 (1968), pp. 55–71.

2) C. E. Bell, "Characterization and Computation of Optimal Policies for Operating an M/G/1 Queuing System with Removable Server," *Operations Research*, Vol. 19 (1971), pp. 208–218.

3) H. Chernoff and A. J. Petkau, "Optimal Control of a Brownian Motion," *SIAM Journal of Applied Mathematics*, Vol. 34 (1978), pp. 717–731.

4) B. T. Doshi, "Controlled One Dimensional Diffusions with Switching Costs—Average Cost Criterion," *Stochastic Processes and Their Applications*, Vol. 8 (1978), pp. 211–227.

5) M. J. Faddy, "Optimal Control of Finite Dams: Continuous Output Procedure," *Advances in Applied Probability*, Vol. 6 (1974), pp. 689–710.

6) W. Feller, "Diffusion Processes in One Dimension," *Transactions of the American Mathematical Society*, Vol. 77 (1954), pp. 1–31.

7) J. M. Harrison and A. J. Taylor, "Optimal Control of a Brownian Storage System," *Stochastic Processes and Their Applications*, Vol. 6 (1978), pp. 179–194.

8) D. P. Heyman, "Optimal Operating Policies for M/G/1 Queuing Systems," *Operations Research*, Vol. 16 (1968), pp. 362–382.

9) D. P. Heyman, "A Diffusion Model Approximation for the GI/G/1 Queue in Heavy Traffic," *The Bell System Technical Journal*, Vol. 54 (1975), pp. 1637–1646.

10) T. Kimura, K. Ohno and H. Mine, "Diffusion Approximation for GI/G/1 Queueing Systems with Finite Capacity: I—The First Overflow Time," *Journal of the Operations Research Society of Japan*, Vol. 22 (1979), pp. 41–68.

11) P. Mandl, *Analytical Treatment of One-dimensional Markov Processes*, Springer-Verlag, New York, 1968.

12) J. H. Rath, "The Optimal Policy for a Controlled Brownian Motion Process," *SIAM Journal of Applied Mathematics*, Vol. 32 (1977), pp. 115–125.

13) S. M. Ross, *Applied Probability Models with Optimization Applications*, Holden-Day, San Francisco, California, 1970.

14) M. J. Sobel, "Optimal Average-Cost Policy for a Queue with Start-Up and Shut-Down Costs," *Operations Research*, Vol. 17 (1969), pp. 145–162.

15) W. Whitt, "A Diffusion Model for a Queue with Removable Server," Technical Report, Department of Administrative Sciences, Yale University, 1973.

16) M. Yadin and P. Naor, "Queueing Systems with a Removable Service Station," *Operational Research Quarterly*, Vol. 14 (1963), pp. 393–405.