

Minimal cost service rate in priority queuing models for emergency cases in hospitals

Nse. S. Udoh^{1*}, Idorenyin A. Etukudo²

¹ Department of Mathematics and Statistics, University of Uyo, Nigeria

² Department of Mathematics and Statistics, Akwa Ibom State University, Nigeria

*Corresponding author E-mail: nseudoh@uniuyo.edu.ng

Abstract

Performance measures and waiting time cost for higher priority patients with severe cases over lower priority patients with stable cases using preemptive priority queuing model were obtained. Also, a total expected waiting time cost per unit time for service and the expected service cost per unit time for priority queuing models: M/M/2: ∞ /NPP and M/M/2: ∞ /PP were respectively formulated and optimized to obtain optimum cost service rate that minimizes the total cost. The results were applied to obtain optimum service rate that minimizes the total cost of providing and waiting for service at the emergency consulting unit of hospital.

Keywords: Priority Queuing Models; Preemptive Model; No Preemptive Model. Waiting Time Cost; Optimum Cost Service Rate; Emergency Case.

1. Introduction

Two priority queuing models; the preemptive priorities (PP) and the Nonpreemptive priorities (NPP) have been considered. In PP, a lower-priority patient with stable case being served can be ejected back into the queue or interrupted whenever a higher-priority patient with serious/critical case, where prompt treatment is vital enters the queuing system. In NPP, a patient being served cannot be ejected back into the queue even if a higher-priority customer enters the queuing system despite the severity of the case.

In order to solve this queuing problem, service facility must be organized so that an optimum balance is obtained between the cost of waiting time and the cost of idle equipment, Gupta and Hira (2009). The cost of waiting-patients generally includes either direct cost of idle doctors and other medical personnel or indirect cost of loss of goodwill due to dissatisfactory service. However, the cost of idle service facilities is the payment to be made for the period.

Hagen et al (2013) examined several queuing models for intensive care units and the effects on waiting times, utilization, return rates, mortalities and number of patients served. Then, Li, et al (2013) designed a hybridized-queuing model for effective packets scheduling in inter-vehicular communication with a view to obtaining the difference between delay-sensitivity and packet length of services. A study of five schemes proved that the nonpreemptive short-packet-first scheme results in the minimal overall delay. The model was superior in terms of performance indices in packet delivery ratio and throughput. Also, Ke, Li and Ni (2012) used the priority queuing model to study the performance of various strategies based on delay sensitivity and packets length. The non-preemptive short-packets-first strategy was proved to result in the minimal overall delay with different strategies for delay-sensitive and non-delay sensitive services. Consequently, an optimal priority-queuing model for the scheduling of multiple internet services was designed based on the above conclusions. Siddarthan, Jones and Johnson (2006) investigated the increased waiting time costs imposed on so-

ciety due to inappropriate use of the emergency department by patients seeking non-preemptive or primary care. It proposed a simple economic model to illustrate the effect of this misuse at a public or not-for-profit hospital. The result showed that non-emergency patients contribute to lengthy delays in the emergency department for all classes of patients. It therefore proposed a priority queuing model to reduce average waiting times.

Several other works on the analysis of emergency waiting time and queuing systems with priority service discipline abound in the literature; see, for example, Udoh (2010), Bedford and Zeepongsekkul (2003), Blake and Carter (1996), among others. The objective of this work is to formulate and optimize the total cost of waiting for service and service cost per unit time to determine the particular level of service which minimizes the total cost of providing service and waiting for that service. We also seek to obtain the associated waiting time costs for the higher priority patients and the lower priority patients.

2. Problem formulation

The waiting time and its associated losses can be decreased by increasing the investment in service facilities such as Doctors and other medical facilities. It is desirable to obtain the minimum sum of these two cost: cost of investment and operation, and cost due to waiting by patients for service. The optimum balance of costs can be obtained by scheduling the flow of patients requiring service and/or providing proper number of doctors and facilities. Both service facilities and the flow of the patients can be controlled by effective schedule of consulting time as well as provide necessary number of doctors and facilities to minimize the overall cost. According to Taha (2007), the estimated waiting time and idle time cost can be obtained as follows;

Let C_w = expected waiting cost/patient/unit time

L_s = expected number of patients in the system

C_t = cost of treating one patient/unit time (that is, cost per unit time Of having a doctor available

μ = service rate or average number of patients treated/unit time

λ = average number of patients' arriving/unit time in the queue

The expected waiting cost per unit time;

$$W_c = C_w L_s = \frac{\lambda}{\mu - \lambda} \tag{1}$$

The expected service cost per unit time;

$$W_s = C_s \mu S_c \tag{2}$$

In addition, total cost

$$C = C_w \frac{\lambda}{\mu - \lambda} + C_s \mu \tag{3}$$

The minimization condition is $\frac{\partial C}{\partial \mu} = 0$

$$\therefore \hat{\mu} = \lambda \pm \sqrt{\frac{\lambda C_w}{C_s}} \tag{4}$$

where $\hat{\mu}$ is the minimum cost service rate.

3. Performance measures and estimated waiting time cost for M/M/2: ∞ /FCFS (in the highest priority class)

Notations

Let P_n = probability of exactly n customers in the queuing system

L = expected number of customers in queuing system

L_q = expected queue length (excluding customers being served

W = waiting time in the system (includes service time) for each individual customer; $W = E(w)$

W_q = waiting time in queue (excludes service time) for each individual customer; $W_q = E(w_q)$

N_t = number of servers in the queuing system at time t ($t \geq 0$)

$P_n(t)$ = probability of exactly n customers in queuing system at time t

S ($=2$) = number of servers (parallel service channels) in the queuing system

λ_n = mean arrival rate (expected number of arrival per unit time of new customers completing service per unit time when n customers are in system)

$\rho = \lambda/\mu$ = traffic intensity (utilization factor for the service facility)

C_q = expected waiting cost in queue

C_s = expected waiting cost in the system

Assumed $\lambda = \lambda_n$, then in a steady-state queuing process $L = \lambda W$ From little's law, Hillier and Lieberman (1995);

$$L_q = \lambda W_q \text{ and } W_q = W - 1/\mu \tag{5}$$

$$\therefore W_q = E(w) = \frac{\lambda}{\mu(\mu - \lambda)}$$

$$L_q = \frac{\lambda^2}{\mu(\mu - \lambda)} \tag{6}$$

When $s > 1$
The utilization factor is

$$U_k = \begin{cases} \left(\frac{\lambda}{\mu}\right)^n & \text{for } n = 1, 2, \dots, s \\ \frac{\left(\frac{\lambda}{\mu}\right)^s}{s!} \left(\frac{\lambda}{s\mu}\right)^{n-s} = \frac{\left(\frac{\lambda}{\mu}\right)^n}{s! s^{n-1}} & \text{for } n = s, s_{n+1} \end{cases}$$

If $N \rightarrow \infty$ and $\lambda < s\mu$, so that $\rho = (\lambda/s\mu) < 1$, then

$$P_0 = \frac{1}{1 + \sum_{n=1}^{s-1} \frac{\left(\frac{\lambda}{\mu}\right)^n}{n!} + \frac{\left(\frac{\lambda}{\mu}\right)^s}{s!} \sum_{n=1}^{\infty} \left(\frac{\lambda}{s\mu}\right)^{n-s}}$$

$$= \frac{1}{\sum_{n=0}^{s-1} \frac{\left(\frac{\lambda}{\mu}\right)^n}{n!} + \frac{\left(\frac{\lambda}{\mu}\right)^s}{s!} \frac{1}{\left(1 - \frac{\lambda}{s\mu}\right)}}$$

$$= \left[\sum_{n=0}^{s-1} \frac{1}{n!} \left(\frac{\lambda}{\mu}\right)^n + \frac{1}{s!} \left(\frac{\lambda}{\mu}\right)^s \frac{s\mu}{s\mu - \lambda} \right]^{-1} \tag{7}$$

and $P_n =$

$$\left\{ \begin{array}{l} \frac{\left(\frac{\lambda}{\mu}\right)^n}{n!} P_0 \text{ if } 0 \leq n \leq s \\ \frac{\left(\frac{\lambda}{\mu}\right)^n}{s! s^{n-s}} P_0 \text{ if } n \geq s \end{array} \right\} \dots \tag{9}$$

Furthermore, $L_q = \sum_{n=s}^{\infty} (n-s)P_n$

$$= \sum_{j=0}^{\infty} j P_{s+j} = \sum_{j=0}^{\infty} j \frac{\left(\frac{\lambda}{\mu}\right)^{s+j}}{s!} \rho^j P_0 = P_0 \frac{\left(\frac{\lambda}{\mu}\right)^s}{s!} \rho \frac{d}{d\rho} \left(\sum_{j=0}^{\infty} \rho^j \right)$$

$$= P_0 \frac{\left(\frac{\lambda}{\mu}\right)^s}{s!} \rho \frac{d\left(\frac{1}{1-\rho}\right)}{d\rho} = P_0 \frac{\left(\frac{\lambda}{\mu}\right)^s \rho}{s!(1-\rho)^2}$$

$$W_q = \frac{L_q}{\lambda} \tag{9}$$

$$\frac{\partial C_{NPP}}{\partial \mu} = \frac{PA_s - QA_s^3 \mu^2}{(Q(s\mu)^2 - Hs\mu + P)^2} - \frac{A}{\mu^2} + C_i = 0$$

$$\frac{(PA_s - QA_s^3 \mu^2) \mu^2 - A((Qs^2 \mu^2 - Hs\mu + P)^2) + C_i((Qs^2 \mu^2 - Hs\mu + P)^2) \mu^2}{(Qs^2 \mu^2 - Hs\mu + P)^2 \mu^2} = 0$$

$$(PA_s - QA_s^3 \mu^2) \mu^2 - A((Qs^2 \mu^2 - Hs\mu + P)^2) + C_i((Qs^2 \mu^2 - Hs\mu + P)^2) \mu^2 = 0$$

$$PA_s \mu^2 - QA_s^3 \mu^4 - A(Qs^2 \mu^2 - Hs\mu + P)^2 + 2PQs^2 \mu^2 - 2HPs\mu + H^2 s^2 \mu^2 + P^2) + C_i(Qs^2 \mu^2 - Hs\mu + P)^2 \mu^2 = 0$$

$$(Qs^2 C_i) \mu^6 - (2HQs^3 C_i) \mu^5 + (2PQs^2 C_i - AQs^3 - AQ^2 s^4 + H^2 s^2 + P^2 C_i) \mu^4 + (2AHPs - 2HPs C_i) \mu^3 + (APs - 2APQs^2 - AH^2 s^2) \mu^2 + (2AHPs) \mu - AP^2 = 0$$

The resulting polynomial in (25) is a function of μ which can be solved for optimum service rate μ^* that minimizes the total cost of providing service and waiting for that service in a NPP queuing model by substituting the values of the constants A, H, Q, P, s, and C_i for a given queuing process.

3.2. Waiting time measure and waiting time cost for preemptive priority model: MM/2: ∞ /PP

If preemption is allowed under the condition in NPP, the total expected waiting time in the system (including the total service time) will change to:

$$w_{sk} = \frac{1/\mu}{B_{k-1} B_k} \text{ for } k = 1, 2, \dots, N \text{ and } s = 1$$

$L_k = \lambda_k w_{sk}$ for $k = 1, 2 \dots N$.
 When $s > 1$, w_k can be calculated by an iterative procedure.
 To determine the expected waiting time in the queue (excluding service time) for priority class k, we have;

$$W_q = W_k - \frac{1}{\mu}$$

Also, if $s > 1$, W_k can be calculated by an iterative procedure. The associated expected waiting time cost is; $w_q C_w$.

3.2.1. Minimum cost service rate (μ^*) for PP queuing model

The waiting time cost is;

$$w_c = C_w L_s = C_w \left(\lambda \left(W_q + \frac{1}{\mu} \right) \right) = C_w \lambda \left[\left(\frac{\lambda}{\mu(\mu - \lambda)} \right) + \frac{1}{\mu} \right]$$

where $W_q = w_{sk} - \frac{1}{\mu}$
 Note that preemption does not affect the service process in any way because of lack of memory property of exponential distribution.
 Hence, the expected total service time for any customer is still $\frac{1}{\mu}$.
 and the expected service cost S_c per unit time is; $S_c = C_i \mu$
 The total cost

$$\frac{\partial C_{PP}}{\partial \mu} = C_w \left[\frac{\lambda}{\mu^2 - \mu\lambda} - \frac{(\lambda^2 + \lambda\mu - \lambda^2)(2\mu - \lambda)}{(\mu^2 - \mu\lambda)^2} \right] + C_i = 0$$

$$C_{PP} = w_c + S_c = C_w \left(\frac{\lambda^2}{\mu(\mu - \lambda)} + \frac{\lambda}{\mu} \right) + C_i \mu = C_w \left(\frac{\lambda^2 + \lambda\mu - \lambda^2}{\mu(\mu - \lambda)} \right) + C_i \mu$$

Similarly, from theorem 1:

$$C_{PP} = C_w [\lambda(\mu^2 - \mu\lambda) - (\lambda^2 + \lambda\mu - \lambda^2)] + C_i C_w (\mu^2 - \mu\lambda)^2 = 0$$

$$C_i \mu^4 - (2\lambda C_i) \mu^3 + (\lambda^2 C_i + C_w \lambda) \mu^2 - 2C_w \lambda \mu = 0$$

Also, the resulting polynomial equation in (31) can be solved for real values of μ^* as the optimum service rate which minimizes the total cost of providing service and waiting for service in PP queuing model by substituting the constants C_i and λ .

4. Application

We consider queuing data on emergency consulting unit of University of Uyo Teaching Hospital in Nigeria. The following parameters were obtained thereof; $\lambda = 2.76$; $s = 2$, $C_i = 17.36$; $r = 1.38$ and $C_w = 4.37$. The associated cost component was obtained as the mean of the salaries of patients/patients’s benefactors. By substituting these parameters in the results in (25) and (31) we obtained the respective service rate equations of the non-preemptive and preemptive priority queuing models as follows;

$$N_{pp} : 3,402.56\mu^6 - 4,928\mu^5 - 1,656.47\mu^4 + 2,784.31\mu^3 - 815.21\mu^2 + 445.09\mu - 39.95 = 0$$

$$P_p : 17.36\mu^4 - 1.736\mu^3 + 0.2619\mu^2 - 0.437\mu = 0$$

The solution of (32) yields the following rational values as roots of μ^* for the Nonpreemptive priority queuing model; 1.4447, -0.8607, 0.7011 and 0.1029. The optimum cost service rate which yield the minimum service rate function that minimizes the waiting time cost of patients and the service time cost of doctors is $\mu^* = 0.7011$.
 Similarly, the solution of (33) yields a rational value of $\mu^* = -0.2786 \approx 0$ as the optimum cost service rate that yields a minimum service rate function for the Preemptive priority queuing model that minimizes the waiting time cost of patients and the service time cost of doctors.

5. Conclusion

The total expected waiting time cost per unit time and the expected service time cost per unit time for priority queuing models: M/M/2: ∞ /NPP and M/M/2: ∞ /PP were formulated and optimized to obtain optimum service rate that both minimizes the expected waiting time cost and the expected service time cost per unit time for the models. The resulting polynomial equations were solved to obtain real values of μ as the optimum cost service rate (μ^*) that minimizes the total cost of time of providing service by doctors and waiting for service by patients at the emergency unit in hospital. The optimal value of $\mu^* = 0.7011$ provide the optimum service time that guarantees minimum waiting time cost and service time cost for the non-preemptive priority model while $\mu^* = -0.2786 \approx 0$ provides an optimum service time that guarantees a minimum waiting time cost and service time cost using the Preemptive priority model. This, however suggest an overstretched of available medical facilities.

References

- [1] Bedford, A and Zeepongsekul, P. On a Dual Queuing System with Preemptive Priority Service Discipline. European Journal of Operational Research, 2003.
- [2] Blake, J. T. and Carter, M. W. An Analysis of Emergency Room Waiting Time Issues via Computer simulation. Information System and Operational Research (INFORM) 34(4):263-273, 1996. <https://doi.org/10.1080/03155986.1996.11732308>.
- [3] Hagen, M. S., Jopling, J. K., Buchmen, T. G and Lee, E. K. Priority Queuing Models for Hospital Intensive Care Units and Impact to Severe Case Patients. AMIA Annual Symposium Proceedings; 841-850, 2013.
- [4] Hillier, S., Frederick and Lieberman, J. Gerald. Introduction to Operations Research. Sixth edition McGraw-Hill, Inc. Singapore, 1995.

- [5] Ke, P., Li, Y. X. and Ni, F. Multiple Services Scheduling with Priority-Queuing Model. IJCIT.org or ijcit.uap-bd.edu (Online) 3(1), 2012.
- [6] Li et al. A Priority-Queuing Model for Heterogeneous Traffic Scheduling in Inter-Vehicular Communication. *Information Technology Journal* 12: 419-423, 2013. <https://doi.org/10.3923/ijtj.2013.419.423>.
- [7] Udoh, N. S. Priority Queuing Models for Emergency Cases in the Hospitals. *Journal of National Association of Mathematical Physics*, 16(2): 191-196, 2010.
- [8] Siddarthan, K., Jones, W. J. and Johnson, J. A. A priority Queuing Model to Reduce Waiting Times in Emergency Care. *International Journal of Healthcare Quality Assurance*, 9(5): 10-16, 2006. <https://doi.org/10.1108/09526869610124993>.
- [9] Taha, H. A. *Operations Research: An Introduction*, 7th edition. Pearson Education Inc., India, 2007.