

On Application of Queuing Models to Customers Management in Banking System

Eze, Everestus Obinwanne¹, Odunukwe, Adaora Darlingtina

Department of Mathematics and Statistics,
Caritas University, Amorji-Nike, Enugu state, Nigeria.

Abstract: Queue is a common sight in banks these days especially on Mondays and on Fridays. Hence queuing theory which is the mathematical study of waiting lines or queue is suitable to be applied in the banking sector since it is associated with queue and waiting line where customers who cannot be served immediately have to queue(wait) for service. The aim of this paper is to determine the average time customers spend on queue and the actual time of service delivery, thereby examining the impact of time wasting and cost associated with it. We used the Markovian birth and death process to analyze the queuing model, which is the Multiple servers, single queue, (M/M/S) queuing model to analyze the data collected by observation from a bank and from the results obtained, the arrival rate is 0.1207 and the service rate is 0.156, the probability that the servers are idle is 0.44 which shows that the servers will be 44% idle and 56% busy, the expected number in the waiting line is 0.1361, the expected number in the system is 0.9098. The expected waiting time in the queue is 1.276 and the expected total time lost waiting in one day is 3.2664 hours, the average cost per day for waiting is ₦65.328 and from the calculation of the comparing solutions, the average cost per day from waiting is ₦7.966 which means that there had been a saving in the expected cost of ₦65.328 - ₦7.966 = ₦57.362. This means that with three servers, the average cost from waiting is reduced. Hence we concluded that the aim and objectives of this paper was achieved.

Keywords: Queue, Service Pattern, Markovian Birth and Death Process, Bank, Poisson Distribution.

I. INTRODUCTION

Queue is a common sight in banks especially on Mondays and Fridays. The word queue comes via French and the latin cauda meaning "tail". Customers waiting in line to receive services in any service system are inevitable and that is why queue management has been where the manager faces huge challenge.

Hence, queuing theory is suitable to be applied in the banking system. Since it is associated with queue or waiting line where customers who cannot be served immediately have to queue (wait) for service for a long time and time being a resource ought to be managed effectively and efficiently because time is money. Queuing theory has become one of the most important, valuable and arguable one of the most universally used tool by an operational researcher. It has applications in diverse fields including telecommunications, traffic engineering, computing and design of factories, shops, offices, banks and hospitals.

Queuing theory can also be applied to a variety of operational situations where it is not possible to accurately predict the arrival rate (or time) of customers and service rate (or time) of service facility of facilities.

Some of the analysis that can be derived using queuing theory include the expected waiting time in the queue, the average waiting time in the system, the expected queue length, the expected number of customers served at one time, the probability of balking customers, as well as the probability of the system to be in certain states, such as empty or full. [Patel, R. et al 2012]

Queuing models are used to represent the various types of queuing systems that arise in practice, the models enable in finding an appropriate balance between the cost of service and the amount of waiting. [A. Nafees, 2007]

Queuing models provide the analyst with a powerful tool for designing and evaluating the performance of queuing systems. [Bank, C. et al, 2001]

A queuing model of a system is an abstract representation whose purpose is to isolate those factors that relate to the system's ability to meet service demands whose occurrences and durations are random. [J. Sztrik, 2010]

¹ Corresponding Author: obinwanneeze@gmail.com

Any system in which arrivals place demands upon a finite capacity resource maybe termed as queuing systems, if the arrival times of these demands are unpredictable, or if the size of these demands is unpredictable, then conflicts for the use of the resource will arise and queues of waiting customers will form and the length of these queue depend on two aspects of the flow pattern: First, they depend on the average rate. Secondly, they depend on the statistical fluctuations of this rate. [Klenrock, L. 1975]

In 1909, the first study of queuing theory was done by a Danish mathematician, A.K. Erlang which resulted into the worldwide acclaimed Erlang telephone model. He examined the telephone network system and tried to determine the effect of fluctuating service demands on calls on utilization of automatic dial equipment.

This study is required to investigate the expected waiting time of customers and the actual waiting time in banks, where the gap between the actual and expected waiting time can be analyzed to know how to improve on the efficiency and effectiveness of their bank. Such problems are:

- How poor service facilities has affected the overall bank performance.
- How poor service pattern affects queue discipline.
- How service facilities has affected the time of customers
- How poor service delivery impacts on time.
- How poor service's delivery has affected customers behaviour.

The aim of this study is to determine the amount of average time customers spend on a queue and actual time of service delivery. Therefore the objectives of this study are as follows:-

- To examine the impact of time wasting on the weak performance.
- To improve on the efficiency and effectiveness of their operations.
- To help bank mangers improve customer's satisfaction through queue management.
- To improve on time management.

This paper when completed will be significant to many people and organisations especially banks in Nigeria. First of all, it will add to the literature on queuing theory and management which will be accessed by lecturers and scholars.

Most importantly, bank managers will benefit a lot from this study as they will apply this theory in their various banks, thereby reducing the amount of time spent on queues which might lead to customer's satisfaction and improve on their overall efficiency and effectiveness.

In terms of the analysis of queuing situations, the types of questions in which we are interested in are typically concerned with measures of system performance which includes:

- To what extent does the service time differ from the actual time that customers have to wait before being served?
- To what extent does poor service pattern affect queue discipline?
- To what extent do service facilities affect customer's service?
- To what extent does the average service time affect the overall performance of the bank?
- How does poor service delivery affect customer's behaviour?

In this paper, we will be studying one bank that adopts the M/M/S queuing model. The queue discipline is first in first out (FIFO) and the arrival is strictly random having Poisson distributed arrival times and exponentially distributed service times.

1.1. Queuing Models and Kendall's Notation

In most cases, queuing models can be characterized by the following factors:

- Arrival time distribution: Inter-arrival times most commonly fall into one of the following distribution patterns a Poisson distribution, a Deterministic distribution, or a General distribution. How-ever, inter-arrival times are most often assumed to be independent and memory less, which is the attributes of a Poisson distribution.

Service time distribution: The service time distribution can be constant, exponential, hyper exponential, or general. The service time is independent of the inter-arrival time.

- Number of server: The queuing calculations change depends on whether there is a single server or multiple servers for the queue. A single server queue has one server for the queue. This is the situation normally found in a grocery store where there is a line for each cashier. A multiple server queue corresponds to the situation in a bank in which a single line waits for the first of several tellers to become available.
- Queue Length: The queue in a system can be modelled as having infinite or finite queue length.
- System capacity: The maximum number of customers in a system can be from 1 up to infinity. This includes the customers waiting in the queue.
- Queuing discipline : There are several possibilities in terms of the sequence of customers to be served such as FIFO (First In First Out, i.e. in order of arrival), random order, LIFO (Last In First Out, i.e. the last one to come will be the first to be served), or priorities.

Kendall, in 1953, proposed a notation system to represent the six characteristics discussed above. The notation of a queue is written as: A/B/P/Q/R/Z where A, B, P, Q, R and Z describes the queuing system properties.

- A describes the distribution type of the inter arrival.
- B describes the distribution type of the service times.
- P describes the number of servers in the system.
- Q describes the maximum length of the queue.
- R describes the size of the system population.
- Z describes the queuing discipline.

1.2. Introducing Notations and Their Terminology

State of the system = number of customers in queuing system.

λ = arrival rate

μ =service rate.

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

μ_n = mean service rate for overall system (expected number of servers completing service per unit time) when n customers are in system.

$N(t)$ = number of customers in queuing system at time t.

$P_n(t)$ = probability that exactly n customers are in queuing system at time t.

$P_0(t)$ = probability that there are no customers in queuing system in time t.

S= number of servers.

M/M/S= inter arrival time and inter departure times are exponentially distributed in a queuing system with S servers.

A birth and death process is one that is appropriate for modelling changes in the size of a population. Thus, most queuing models especially the single ones can be analyzed as birth and death processes. In the context of queuing theory, we define the following terms:-

- Birth refers to the arrival or entrance of a new customer into the queuing system.
- Death refers to the departure of a served customer.
- The state of the system at time t, (t>0) is given by N(t)

Thus, the birth and death process describes probabilistically how N (t) changes as t changes.

Generally speaking, the birth and death process states that individual birth and death occurs randomly, where their mean occurrence rates depends only upon the current state of the system.

1.3. Multiple Servers Queue with M|M|S Model

An M|M|S system is a queuing process having Poisson arrival patterns, S server, with S independent, identically distributed, exponential service times (which does not depend on the state of the system); infinite capacity, and a FIFO queue discipline. The arrival pattern being stated independent, $\lambda_n = \lambda$ for all n. The service times associated with each server are also independent, but since the number of servers that actually attend to customers (i.e. are not idle) does depend on the number of customers in the system, the effective time it takes the system to process customer through the service time facility is state dependent. In particular, if $\frac{1}{\mu}$ is the mean service time for one server to handle one customer, then the mean rate of service completion when there are customers in the system is

$$\lambda_n = \lambda, \mu_n = \begin{cases} n\mu & \text{if } n = 0, 1, \dots, c. \quad \text{i.e. } n \leq c \\ c\mu & \text{if } n = c + 1, c + 2. \quad \text{i.e. } n \geq c \end{cases} \quad (1.3.1)$$

The probability of zero customers in the system (P_0) and the probability of n customer in the system (P_n) are given by:

$$P_0 = \left\{ \frac{\left(\frac{\lambda}{\mu}\right)^c}{c! \left[1 - \frac{\lambda/\mu}{c}\right]} + \frac{(\lambda/\mu)^1}{1!} + \frac{(\lambda/\mu)^2}{2!} + \dots + \frac{(\lambda/\mu)^{c-1}}{(c-1)!} \right\}^{-1} \quad (1.3.2)$$

$$P_n = P_0 \frac{(\lambda/\mu)^n}{n!} \quad \text{if } n \leq c \quad (1.3.3)$$

$$P_n = P_0 \frac{(\lambda/\mu)^n}{c! c^{n-c}} \quad \text{if } n > c \quad (1.3.4)$$

The capacity utilization in this system is $\lambda/c\mu$

We can use the above equation of $\lambda/c\mu < 1$

If $\frac{\lambda}{c\mu} > 1$, then the waiting line grows larger and larger i.e. becomes infinite if the process runs long enough.

when $C = 1$ (there is one service facility), equations (1.3.3) and (1.3.4) reduces to

$$P_n = \left(\frac{\lambda}{\mu}\right)^n P_0 \quad (1.3.5)$$

From equations (1.3.3) and (1.3.4), we have

$$P_n = P_0 \frac{\left(\frac{\lambda}{\mu}\right)^n}{n!} \quad \text{if } n \leq C$$

But n can only take on values of 0 or 1 if $n \leq C = 1$. Thus

$$P_n = P_0 \left(\frac{\lambda}{\mu}\right)^n$$

If $C = 1$, equation 4 also reduces to equation 5

With C service facilities, the average number of customers in the queue is

$$N_q = \frac{(\lambda/\mu)^{c+1} P_0}{c.c! \left[1 - \frac{\lambda/\mu}{c}\right]^2} \quad (1.3.6)$$

The average number in the system (waiting plus service) is

$$N_s = N_q + \frac{\lambda}{\mu} \quad (1.3.7)$$

The expected waiting time in the queue for an arrival is

$$T_q = \frac{N_q}{\lambda} \quad (1.3.8)$$

The expected total time spent in the system (waiting plus service) is

$$T = \frac{N_s}{\lambda} \quad (1.3.9)$$

II. METHODOLOGY

The method adopted in this paper is the Birth and Death process which is a special class of Markov process.

III. MAIN RESULT

The methodology can now be illustrated using the data below: if the data collected by observation in a bank shows that the system capacity is 238 customers with 2 servers, inter-arrival time for 238 customers is 1972 minutes and the time taken by 239 customers to be served is 1529 minutes.

We first of all estimate the parameters,

N= 238 customers

T= 1972 minutes

S= 1529 minutes

C= 2

Where N= system capacity

T= time interval

S= time taken by 239 customers to be served

C= number of servers.

Now the

$$\text{Arrival rate, } \lambda = \frac{N}{T} = \frac{238}{1972} = 0.1207$$

$$\text{Service rate, } \mu = \frac{N}{S} = \frac{238}{1529} = 0.156$$

$$\text{Traffic intensity } (\rho) = \frac{\lambda}{\mu} = \frac{0.1207}{0.156} = 0.7737$$

This implies that

$$\mu = 0.1207$$

$$\lambda = 0.156$$

$$C = 2$$

$$\frac{\lambda}{\mu} = 0.7737$$

We now calculate the probability values

- Probability that the servers are idle is equ (1.3.2)

$$P_0 = \frac{1}{\left(\frac{\lambda}{\mu}\right)^c + 1 + \frac{(\lambda/\mu)^1}{1!} + \dots + \frac{(\lambda/\mu)^{c-1}}{(c-1)!}$$

$$P_0 = \frac{1}{\frac{0.7737}{2! \left[1 - \frac{0.7737}{2}\right]} + 1 + \frac{0.7737}{1}}$$

$$= \frac{1}{0.4881 + 1.7737} = \frac{1}{2.2618} = 0.4421$$

- Probability of n customers in the system is from equ. (1.3.4)

$$P_n = \frac{\left(\frac{\lambda}{\mu}\right)^n}{c! c^{n-c}} = \frac{(0.7737)^n 0.4421}{2! \times 2^{n-2}} = \frac{(0.7737)^n 0.4421}{2^{n-1}}$$

Now we calculate the Queue measure

- The expected number in the waiting line, from equ (1.3.6) we have that

$$N_q = \frac{\left(\frac{\lambda}{\mu}\right)^{c+1}}{c c! \left[1 - \frac{\lambda}{c\mu}\right]^2} P_0$$

$$= \frac{0.7737^3 (0.4421)}{2 \times 2! \left[1 - \frac{0.7737}{2}\right]^2} = \frac{(0.4631) \times 0.4421}{4 \times 0.3760} = 0.1361$$

- The expected number in the system (waiting plus in service) is by equ. (1.3.7)

$$N_s = N_q + \frac{\lambda}{\mu}$$

$$= 0.1361 + 0.7737 = 0.9098$$

- The expected waiting time in the queue is by equation (1.3.8)

$$T_q = \frac{N_q}{\lambda} = \frac{0.1361}{0.1207} = 1.1276$$

- The expected total time lost in one day

$$T_L = \lambda \times 24 \text{ hours} \times T_q \text{ (hint 1 day = 24 hours)}$$

$$= 0.1207 \times 24 \times 1.1276 = 3.2664 \text{ hours.}$$

Assuming a cost is associated with this hour (i.e. ₹20 for each hour lost by customer waiting).

The average cost per day from waiting is:

$$= 3.2664 \times ₹20 = ₹65.328.$$

Now, we want to find out if increasing the number of servers can help to reduce the amount of time spent on queue and the hence minimize the cost incurred by waiting. Hence we compare solutions:

Let $\lambda = 0.1207$

$$\mu = 0.156$$

$$\frac{\lambda}{\mu} = 0.7737 \text{ which is same in the data above but let } C = 3$$

Now we have

$$P_0 = \frac{1}{\frac{(0.7737)^3}{3! \times \left[1 - \frac{0.7737}{3}\right]} + 1 + \frac{(0.7737)^1}{1!} + \frac{(0.7737)^2}{2!}}$$

$$= \frac{1}{0.1040 + 1 + 0.7737 + 0.2993} = \frac{1}{2.177} = 0.4593$$

The expected number in the waiting line is

$$N_q = \frac{(0.7737)^4}{3 \times 3! \left(1 - \frac{0.7737}{3}\right)^2} (0.4593)$$

$$= \frac{0.1646}{9.9128} = 0.0166$$

Expected number in the system is

$$N_s = N_q + \frac{\lambda}{\mu}$$

$$= 0.0166 + 0.7737 = 0.7903$$

Then expected waiting time in the queue is

$$T_q = \frac{N_q}{\lambda}$$

$$= \frac{0.0166}{0.1207} = 0.1375$$

The expected total time lost waiting on one day

$$T_L = \lambda \times 24 \text{ hours} \times T_q$$

$$= 0.1207 \times 24 \times 0.1375 = 0.3983 \text{ hours}$$

The average lost per day from waiting is

$$= 0.3983 \times \text{₦}20 = \text{₦}7.966$$

IV. DISCUSSION AND CONCLUSION

Considering the analytical solution, the capacity of the system under study is 238 customers and the arrival rate is 0.1207 while the service rate is 0.156. This shows that the service rate of the system is greater than the arrival rate, this does not necessarily imply that there is no queue but that queue may not be long. Considering the ratio of the two rates that is arrival rate over the service rate called traffic intensity which is less than one. Probability that the servers are idle is 0.44 which shows that the servers will be 44% idle and 56% busy, while the probability that there are n customers in the system is

$$\frac{(0.4421)(0.7737)^n}{2^{n-1}}$$

The expected number in the waiting line is 0.1361. The expected number in the system is 0.9098. The expected waiting time in the queue is 1.276 and the expected total time lost waiting in one day is 3.2664 hours.

From the foregoing Queue measures, the average cost per day for waiting is ₦65.328 and from the calculation of the comparing solutions, the average cost per day from waiting is ₦7.966. There had been a saving in the expected cost of ₦65.328 - ₦7.966 = ₦57.362. This means that with three servers, the average cost from waiting is reduced.

We now conclude that adding one more server will help reduce time spent on queue which can improve customer's satisfaction. Hence the objective of this paper is achieved.

The advantage of using this single system with multiple servers is that a slow server does not affect the movement of the queue i.e if a server is slow it does not affect the movement of the queue because next customer can go to the next available server instead of waiting for the slow server.

REFERENCES

- [1] A.k Sharma, G.K Sharma (2013): Queuing theory approach with queuing model: A study" ISSN: 231966726 volume 2 issue 2, February 2013 PP: 01-11
- [2] Banks, J. Carson, J.S. Nelson, B.L, Nicol, D.M (2001): Discrete – Event System Simulation, Prentice Hall international series, 3rd edition, P.24-37.
- [3] Dr. Janos Sztrik (2011): Basic Queuing theory. University of Debrecen, Faculty of Informatics.
- [4] Kasumu R.B (1994); Introduction to probability theory; Fatol Ventures Publishers. Lagos Nigeria.
- [5] Kasumu R.B (2000); Introduction to Stochastic Process; Fatol Ventures Publishers Lagos Nigeria.
- [6] Kleinrock, L (1975): Queuing systems. Vol. I. Theory. John Wiley & Sons, New York.
- [7] O. Olasore (1992); Nigeria Banking and Economic Management "Journal of the Nigeria Institute of Bankers" Vol.3 page 8-13. Nigeria.
- [8] Patel, J.J, Rajeshkumar, M.C, Pragnesh, A.P, Makwana, P (2012) "Minimise the waiting time of customer and gain more profit in restaurant using queuing model" International conference on management, humanity and economics(ICMHE 2012) p77-80