# Queuing Theory

A flow of customers from finite/infinite population towards the service facility forms a queue on account of lack of capability to serve them all at one time. In the absence of a perfect balance between the service facilities and the customers, waiting is required either of the service facility or for the customers. The subject of queueing is not directly concerned with optimization, it attempts to explore, understand and compare various queuing situations, thus indirectly achieving approximate optimization.

## Elements of queueing system

1. Input process - This element of a queueing system is concerned with the pattern in which the customers arrive for service. This can be described by the following three factors:

(a) Size of queue - If the total no of customers requiring service are few, then the size is said to be finite. If the potential customers are sufficiently large, then the source is considered infinite. If more than one arrival is allowed to enter the system simultaneously, the input is said to occur in bulk or patches.

(b) Arrival distribution - If the pattern of arrival of customers at the service system is uncertain between successive

arrivals (inter-arrival time), then the arrival pattern is measured by mean arrival time. These are characterised by the probability distribution associated with this random process. Most queueing models assume that the arrival rate follows a Poisson distribution and the inter-arrival time follows an exponential distribution.

(c) <u>Customer's behaviour</u> — A customer may decide to wait no matter how long the queue is or he may decide not to enter the system if the queue is too long. If he decides not to wait, he is said to have '<u>balked</u>'. He may enter the queue but decide to leave after some time. In this case, he is said to have '<u>reneged</u>'. In case there is more than one queue, he may shift between queues. This is called '<u>jockey' for position'</u>.

(d) The input process which does not change with time is called a <u>stationary</u> process. If it is time dependent, it is called transient.

2. <u>Queue Discipline</u> — This is the rule according to which, customers are served after the queue is formed. The most common dispt discipline is FIFO or FCFS. Other disciplines include LIFO and SIRO (service in random order). There are

variety of priority schemes in which service is done in preference (priority scheme). Under this case, service is of two types:

(i) pre-emptive — Service is provided to customers of high priority

(ii) Non pre-emptive — Service is provided to customers of low priority

In case of parallel channels. FSR (Faster server first Rule) is adopted.

**Service mechanism** — The service mechanism is concerned with service time that is the time interval from the commencement of service to completion of service. Service facilities are of the following type:

(a) single queue - single server
(b) single queue - several server
(c) several queue - single server
(d) ~~several queue - several server~~
(e) several servers ⎯ ⎯ series channels / parallel channels

**Capacity of system**

The source from which the customers are generated can be finite / infinite

## Operating characteristics

$E(n) = L$ : expected no. of customers in the system.

$E(m) = Lq$ : expected no. of customers in queue ; $m = n-1$

$E(v) = W$ : expected waiting time in the system.

$E(w) = Wq$ : expected waiting time in the queue

Server utilisation factor, $P = \frac{\lambda}{\mu} \neq$ A where

A = average customers arriving per unit time, $\mu$ = average customers completing service per unit time. $P$ is also called traffic intensity.

## Probability distribution in queueing system

It is assumed that customers joining the system arrive in a random manner and follow a Poisson distribution the service times are mostly assumed to be exponentially distributed. It implies the probability of service completion is independent of length of time spent in progress. We assume the arrival and service distributions follow a Poisson Queue, under the following axioms.

Axiom 1 : The no. of arrivals in non-overlapping time interval are statistically independent.

Action 2 : The probability of more than one arrival within a small time interval, $[t, t+\Delta t]$ is negligible.

$$P_0(\Delta t) + P_1(\Delta t) + O(\Delta t) = 1$$

$$P_n(t) = p \ (n \ arrivals \ in \ a \ time \ interval, \ t)$$

Action 3 : The probability that an arrival occurs rate within time $[t, t+\Delta t]$ is

$$P_1(\Delta t) = \lambda \Delta(t) + O(\Delta t)$$

[ derivations — study from book
  (a) distribution of arrival
  (b) distribution of departure ]

$(M/M/1) : (\infty / FIFO)$   // infinite source
$(M/M/1) : (N / FIFO)$   // finite source.

Pg. 592 [21.5 — "Deterministic Queuing system"
XX 2 Distribution of Inter-arrival XX
└ 3 Di "   " departure

21:7  classification of Queuing models

21:9  model 1 — derivation not need
   characteristics of model 1 — formulae

Relationship  X model II X

Pg 608  model III — only eqⁿs.