

Priority queues

Consider an $M/G/1$ queue where the customers are divided into K priority classes, $k = 1, \dots, K$:

- class 1 has the highest priority and class K the lowest priority
- the arrival rates of different classes are $\lambda_1, \dots, \lambda_K$ (Poissonian)
- the expectation and second moment of the service time of different classes: \overline{S}_k ja \overline{S}_k^2 , $k = 1, \dots, K$.

We derive mean results of the Pollaczek-Khinchinin type for this kind of queueing system.

- Priority queues are becoming more important also in telecommunication systems
 - in computer systems (e.g. operating systems) they have been exploited for a long time
- We are focusing on the so called time priority, which defines the service order of different classes
 - by giving a customer a higher priority one can reduce delay and delay variation
- For finite size queues (e.g. a finite buffer), a separate issue arises regarding the control of losses (overflow); this gives rise to the notion of space priority.

Non-preemptive priority

- The service of the customer being served is completed even if customers of higher priority may arrive.
- Each priority class has a separate (logical) queue.
- When the server becomes free, customer from the head of the highest priority non-empty queue enters the server.

Notation

$$\left\{ \begin{array}{l} \bar{N}_q^{(k)} = \text{mean number of waiting class-}k \text{ customers in the queue} \\ \bar{W}_k = \text{mean waiting time of class-}k \text{ customers} \\ \rho_k = \text{the load of class } k, \rho_k = \lambda_k \bar{S}_k \\ \bar{R} = \text{the mean residual service time in the server (upon arrival)} \end{array} \right.$$

The stability condition of the queue:

$$\rho_1 + \cdots + \rho_K < 1$$

If the condition is violated, the queues belonging to priority classes lower than some limit k (i.e. with a higher class index k) will grow without bound.

Non-preemptive priority (continued)

In the same way as in the derivation of the Pollaczek-Khinchin mean results we deduce for the highest priority class 1:

$$\bar{W}_1 = \bar{R} + \bar{S}_1 \bar{N}_q^{(1)} \quad \text{where the latter term represents the average time needed to serve the class-1 customers ahead in the queue}$$

By Little's result we have

$$\bar{N}_q^{(1)} = \lambda_1 \bar{W}_1 \quad \Rightarrow \quad \bar{W}_1 = \bar{R} + \rho_1 \bar{W}_1 \quad \Rightarrow \quad \boxed{\bar{W}_1 = \frac{\bar{R}}{1 - \rho_1}}$$

For priority class 2 we get analogously

$$\bar{W}_2 = \bar{R} + \underbrace{\bar{S}_1 \bar{N}_q^{(1)} + \bar{S}_2 \bar{N}_q^{(2)}}_{\text{time needed to serve class-1 and class-2 customers ahead in the queue}} + \underbrace{\bar{S}_1 \lambda_1 \bar{W}_2}_{\text{time needed to serve those customers in higher classes that arrive during the waiting time of class-2 customer}}$$

By Little's result we again get

$$\bar{N}_q^{(2)} = \lambda_2 \bar{W}_2 \quad \Rightarrow \quad \bar{W}_2 = \bar{R} + \rho_1 \bar{W}_1 + \rho_2 \bar{W}_2 + \rho_1 \bar{W}_2 \quad \Rightarrow \quad \bar{W}_2 = \frac{\bar{R} + \rho_1 \bar{W}_1}{1 - \rho_1 - \rho_2}$$

Non-preemptive priority (continued)

By inserting in the expression for \overline{W}_2 the formula for \overline{W}_1 we get

$$\overline{W}_2 = \frac{\overline{R}}{(1 - \rho_1)(1 - \rho_1 - \rho_2)}$$

By continuing in the same way to lower priority classes (higher values of k) we get the general result

$$\overline{W}_k = \frac{\overline{R}}{(1 - \rho_1 - \dots - \rho_{k-1})(1 - \rho_1 - \dots - \rho_k)}$$

The total time in system of class- k customers is on the average

$$\overline{T}_k = \overline{W}_k + \overline{S}_k$$

The mean residual service time \overline{R} appearing in \overline{W}_k can be derived by the same kind of graphical “triangle trick” as in the case of the Pollaczek-Khinchin mean value formula:

$$\overline{R} = \frac{1}{2} \sum_{k=1}^K \lambda_k \overline{S}_k^2$$

Observations about the non-preemptive priority

- The analysis cannot easily be extended for multiserver systems
 - the residual service time \bar{R} is difficult to determine
 - can, however, be determined if the service time of all the classes has an exponential distribution with the same mean
- The mean waiting time of customer can be controlled by the choice of the priority class
 - if highest priority is given to those customers which have the shortest service times, then the mean waiting time of the whole customer population will be minimized
 - cf. the queue at a copying machine, where usually a person with only one page to copy is given priority over those taking copy of a whole report
 - in the case of two priority classes the mean sojourn time is

$$\bar{T} = \frac{\lambda_1 \bar{T}_1 + \lambda_2 \bar{T}_2}{\lambda_1 + \lambda_2}$$

- one can show that if $\bar{S}_1 < \bar{S}_2$, then \bar{T} is smaller than in the case where the priorities were interchanged (or no priorities were used)
- The waiting time of also the highest priority class 1 depends on the traffic of the lower classes (on the λ_k) because the non-preemptivity means that the lower classes are not completely “invisible” to the higher classes.

Kleinrock's conservation theorem for non-preemptive priority queues

Kleinrock's conservation theorem states

$$\boxed{\sum_{k=1}^K \rho_k \bar{W}_k = \frac{\rho \bar{R}}{1 - \rho}}$$

where $\rho = \rho_1 + \dots + \rho_K$ is the total load of the system.

- The weighted sum of the waiting times can never change no matter how sophisticated the queueing discipline may be.
- Any attempt to modify the queueing discipline so as to reduce one of the \bar{W}_k will force an increase in some of the other \bar{W}_k .

In particular, if all the classes have the same average service time, $\bar{S}_k = s$, one obtains upon division by s

$$\frac{\lambda \bar{R}}{1 - \rho} = \sum_{k=1}^K \lambda_k \bar{W}_k \stackrel{\text{(Little)}}{=} \sum_{k=1}^K \bar{N}_q^{(k)} = \bar{N}_q$$

The average total number of waiting customers \bar{N}_q is a constant independent of the queueing discipline (under the assumption of the same average service time).

Proof of Kleinrock's conservation theorem

The average unfinished work \bar{V} can be divided into two parts

$$\bar{V} = \bar{V}_q + \bar{R}$$

The average unfinished work in the queue (that of the waiting customers) \bar{V}_q can be written with the aid of Little's theorem

$$\bar{V}_q = \sum_k \bar{N}_q^{(k)} \bar{S}_k = \sum_k \lambda_k \bar{W}_k \bar{S}_k = \sum_k \rho_k \bar{W}_k$$

Thus

$$\sum_k \rho_k \bar{W}_k = \bar{V} - \bar{R}$$

\bar{V} is independent of the service order (holds for the whole process $V(t)$). So is \bar{R} since all the customers ultimately go through the server. This shows that $\sum_k \rho_k \bar{W}_k$ is constant.

Without changing the value we can calculate \bar{V} for any service discipline. Let us calculate it for an ordinary FIFO system ($M/G/1$ queue) which consists of a single class.

$$\bar{V} = \sum_k \rho_k \bar{W}_k + \bar{R} = \rho \bar{W} + \bar{R} \stackrel{\text{(PASTA)}}{=} \rho \bar{V} + \bar{R} \quad \Rightarrow \quad \bar{V} = \frac{\bar{R}}{1 - \rho}$$

By substituting in the previous equation one obtains the desired result $\sum_k \rho_k \bar{W}_k = \frac{\rho \bar{R}}{1 - \rho}$

Preemptive resume priority

The service of customer is interrupted when on arrival of customer belonging to a higher class arrives, and will be resumed from the point of interruption, when all the queues of the higher priority classes have been emptied

- In this case the lower priority class customers are completely “invisible” and do not affect in any way the queues of the higher classes.
- In the case of packet queues, the preemptive priority is approximated if the sending of the packets takes place in small pieces, e.g. ATM cells which have (non-preemptive priorities)
 - when a packet of a higher priority arrives, the sending of the cells belonging to the lower priority packet is interrupted and is continued first when all the higher priority packets have been fully sent

Preemptive resume priority (continued)

We wish to calculate the mean sojourn time \bar{T}_k of class- k customers. It is composed of three parts:

1. The customer's own mean service time \bar{S}_k
2. The mean time to serve the customers in classes $1, \dots, k$ ahead in the queue

$$\frac{\bar{R}_k}{1 - \rho_1 - \dots - \rho_k}, \text{ missä } \bar{R}_k = \frac{1}{2} \sum_{i=1}^k \lambda_i \bar{S}_i^2,$$

which is the same as the mean waiting time in an $M/G/1$ queue arising solely from the customers in classes $1, \dots, k$. This is due to the fact that the unfinished work of classes $1, \dots, k$ (not affected by classes $k+1, \dots, K$) does not at all depend on the mutual service order of customers in classes $1, \dots, k$ (since the unfinished work is “anonymous”),

$$E^*[V_k(t)]_{\text{pr. queue}} = E^*[V_k(t)]_{M/G/1} = E[W_k]_{M/G/1}$$

3. The mean time it takes to serve those customers in the higher classes $1, \dots, k-1$ which arrive during the time the class- k customer stays in the system

$$\sum_{i=1}^{k-1} \bar{S}_i \lambda_i \bar{T}_k = \sum_{i=1}^{k-1} \rho_i \bar{T}_k, \quad k > 1 \quad (0, \text{ jos } k = 1)$$

Preemptive resume priority (continued)

By collecting the results together we obtain

$$\bar{T}_k = \bar{S}_k + \frac{\bar{R}_k}{1 - \rho_1 - \cdots - \rho_k} + \left(\sum_{i=1}^{k-1} \rho_i \right) \bar{T}_k$$

$\bar{T}_1 = \frac{(1 - \rho_1)\bar{S}_1 + \bar{R}_1}{1 - \rho_1}$
$\bar{T}_k = \frac{(1 - \rho_1 - \cdots - \rho_k)\bar{S}_k + \bar{R}_k}{(1 - \rho_1 - \cdots - \rho_{k-1})(1 - \rho_1 - \cdots - \rho_k)}$