How do we choose the rates in an Accumulating Priority Queue?

Azaz Bin Sharif*, David Stanford*, Peter Taylor[†] and Ilze Ziedins[‡] *University of Western Ontario, [†]University of Melbourne, [‡]University of Auckland.



CTAS

The Canadian Emergency Department Triage and Acuity Scale

(http://www.calgaryhealthregion.ca/policy/docs/1451/Admission_over-capacity_AppendixA.pdf)

Category	Classification	Access	Performance Level
1	Resuscitation	Immediate	98%
2	Emergency	15 minute	95%
3	Urgent	30 minute	90%
4	Less urgent	60 minute	85%
5	Not urgent	120 minute	80%



The Australasian Triage Scale

(http://www.acem.org.au/media/policies and guidelines/P06 Aust Triage Scale - Nov 2000.pdf)

Category	Access	Performance Level
1	Immediate	100%
2	10 minute	80%
3	30 minute	75%
4	60 minute	70%
5	120 minute	70%

In both cases, there are priority classifications, with associated access targets and proportions of time that the target should be met.

The KPIs are in terms of tails of waiting time distributions.



What would we like to know?

- 1. We would like to know whether it is possible to meet the targets for all customer classes simultaneously.
- 2. If so, we want to propose a queueing discipline that ensures that the targets are met.



The accumulating priority queue

In 1964, Kleinrock proposed a queueing discipline where

- customers accumulate priority at class-dependent rates as linear functions of their time in the queue.
- when the server becomes free, it selects the waiting customer with the highest amount of accumulated priority at that instant, provided that the queue is non-empty.



The accumulating priority queue

- There is a single server.
- Customers of priority *i* arrive according to independent Poisson streams with rate λ_i.
- They accumulate priority at rate b_i where $1 = b_1 > b_2 > \ldots > b_l > 0$.
- When the server becomes free it chooses to serve the customer with the highest current priority, if the system is non-empty.
- Service times are chosen independently from the class-dependent distribution function B_i(t) with Laplace-Stieltjes Transform B^{*}_i(s).



The APQ





Mean waiting Times

For such a queue, Kleinrock developed a recursion for calculating the expected waiting times of customers from each class.

He showed that

$$W_{i} = \frac{[\hat{M}/(1-\rho)] - \sum_{j=i+1}^{l} \rho_{j} W_{j} [1-b_{j}/b_{i}]}{1 - \sum_{j=1}^{i-1} \rho_{j} [1-b_{i}/b_{j}]}$$

where

- *M_i* is the mean service time of type *i* customers,
- $\rho_i = \lambda_i M_i$ with $\rho = \sum_{i=1}^{l} \rho_i$, and
- \hat{M} is the stationary mean residual service time.



Mean waiting Times

The above recursion can be inverted to deliver the b_i that produce desired ratios W_i/W_1 of the mean stationary waiting times (within a range restricted by these ratios for a pure priority system).

So if our performance standards were specified in terms of

- mean waiting times, rather than
- tails of waiting time distributions,

we would have the answer to both of our questions.



Waiting time distributions

Theorem (Stanford, Taylor and Ziedins (2014))

For a single-server APQ, the LST $\tilde{W}^{(i)}_+(s)$ of the waiting time distribution for a class-*i* customer, conditional on it being positive, is given by

$$ilde{W}^{(i)}_+(s) = (1-b_{i+1}/b_i) ilde{W}^{(i)}_{acc}(s) + (b_{i+1}/b_i) ilde{W}^{(i+1)}_+(b_{i+1}s/b_i)$$

where

$$egin{array}{rcl} ilde{W}^{(i)}_{acc}(s) &=& rac{1-
ho}{1-\delta_i}\, ilde{W}^{(i,0)}_{acc}(s) + \, ilde{W}^{(i+1)}_+(b_{i+1}s/b_i)\sum_{j=1}^{l}\,rac{
ho_j b_{i+1}}{b_j(1-\delta_i)}\, ilde{W}^{(i,j)}_{acc}(s) \ &+& \sum_{j=i+1}^{l}rac{
ho_j}{1-\delta_i}\, ilde{W}^{(j)}_+(b_is/b_j) ilde{W}^{(i,j)}_{acc}(s), \end{array}$$

and $\delta_i = \sum_{j=1}^i \rho_j (1 - b_{i+1}/b_j)$.



Choosing the accumulation rates

For a given set of workload parameters $\rho = (\rho_1, \dots, \rho_l)$ and tuning parameters $\mathbf{b} = (b_1, \dots, b_l)$, we can

- use the STZ recursion to evaluate the Laplace transform of the waiting time at any point *s* that we are interested in,
- numerically invert the transform to derive the waiting time distributions W⁽ⁱ⁾(t), and
- test whether the performance standards in terms of tails are met.



Choosing the accumulation rates



Class 1 Waiting Time CDF



Choosing the accumulation rates

So,

- we have a way of testing whether any $\mathbf{b} = (b_1, \dots, b_l)$ is feasible, but
- we do not have a way of defining the feasible region for **b** = (b₁,..., b_l) analytically.

• We also haven't specified any objective function.



David Stanford gave a talk at the 'Queues, Modelling and Markov Chains' Workshop in which he suggested that we should minimise the expected waiting time of those patients that do not meet the performance KPIs.

Category	Classification	Access	Performance Level
1	Resuscitation	Immediate	98%
2	Emergency	15 minute	95%
3	Urgent	30 minute	90%
4	Less urgent	60 minute	85%
5	Not urgent	120 minute	80%

This is a good idea from a 'human' point of view, but



Another possibility is to choose the value of the $\mathbf{b} = (b_1, \dots, b_l)$ that satisfy the constraints for the greatest range of $\rho = (\rho_1, \dots, \rho_l)$.



Feasibility range for increasing ρ with $\lambda_1 = \lambda_2$.









