



Point of queue size change analysis of the PH/PH/k system with heterogeneous servers

Attahiru Sule Alfa^{a,b,*}, Qi-Ming He^c

^a Department of Electrical and Computer Engineering, University of Manitoba, Winnipeg, Manitoba, Canada R3T 5V6

^b Department of Electrical, Electronic and Computer Engineering, University of Pretoria, Pretoria, South Africa

^c Department of Management Sciences, University of Waterloo, Waterloo, Ontario, Canada N2L 3G1



ARTICLE INFO

Article history:

Received 2 March 2017

Received in revised form 8 September 2017

Accepted 14 September 2017

Available online 22 September 2017

Keywords:

PH/PH/k system

Point of queue size change analysis

Continuous time and discrete time

ABSTRACT

By observing the points only when the queue size changes, we study the PH/PH/k system with heterogeneous servers. We develop Markov chain set ups that are more efficient than studying the systems at arbitrary times, by reducing the sizes of matrices that need to be computed. Specifically we present procedures for constructing the associated Markov chains so one may use matrix-analytic methods for their analysis. This work is carried out for both the continuous and discrete time cases.

© 2017 Elsevier B.V. All rights reserved.

1. Introduction

Multi-server queueing systems are natural systems that occur in real life. An example is in telecommunication systems where the servers are communication channels which are usually not identical. Even though such queueing systems can be set up and analyzed by observing them at arbitrary times this leads to huge matrices which are often inefficient to compute and use in the system analysis. When the servers are identical the method of studying the system by considering the number of servers in each phase makes it easier to analyze (see [3,5,9,10]). However when the servers are not identical the methods in those papers cannot be used. Trying to study the systems by the traditional approaches of studying the systems at arbitrary time points would lead to huge block matrices in the associated Markov chains. If, however, we study the systems at points of queue size changes only, a form of embedded system, then we can reduce the block sizes of the associated matrices. The performance measures at arrival (or departure) time points can be obtained from that of the points of queue size changes. Nonetheless, constructing the Markov chains and obtaining the block matrices is a challenging procedure. In this paper we show how to construct these matrices for both the continuous time and discrete time cases. Latouche and Ramaswami [6] presented a method for analyzing the continuous time PH/PH/1 system at points of queue size changes. Their work

is a special case with a single server, and they showed that the method helps to reduce the matrix needed for computation to be of size $(m + n) \times (m + n)$, from $mn \times mn$, where n and m are the dimensions associated with the PH-distributions for the interarrival and service times, respectively. The discrete time case of PH/PH/1 was presented in [1].

2. The continuous time model

We consider the continuous time PH/PH/k queue in this section. In Section 2.1, we define the PH/PH/k queue and an embedded Markov chain at the points of events. In Section 2.2, we construct the block matrices in the transition matrix of the embedded Markov chain. A brief analysis of the queue length and waiting time is presented in Section 2.3.

2.1. PH/PH/k system at points of events

Let the arrival process be PH with representation (β_0, S_0) of order m_0 . There are k servers and the service time of server r is of the PH type with representation (β_r, S_r) of order m_r , $r = 1, 2, \dots, k$. If an arriving customer finds multiple servers available, the customer chooses the server with the smallest index for service. We note that models with other Markovian rules for server selection can be treated with minor modifications of the analysis in this paper.

Let $q(t)$ be the number of customers in the system at time t , $I_0(t)$ be the phase of the PH arrival process at time t , and $I_r(t)$ be the phase of the service process of server r at time t , for $r = 1, 2, \dots, k$. We note that $I_r(t) \in \{1, 2, \dots, m_r\}$, for $r = 0, 1, 2, \dots, k$. For the

* Corresponding author at: Department of Electrical and Computer Engineering, University of Manitoba, Winnipeg, Manitoba, Canada R3T 5V6.

E-mail addresses: attahiru.alfa@umanitoba.ca (A.S. Alfa), q7he@uwaterloo.ca (Q.-M. He).

- If the last event is the completion of a service, we have

$$B_{*0}^{(\xi_1, \dots, \xi_m)} = - \left(S_0 \oplus \bigoplus_{j=1}^m S_{\xi_j} \right)^{-1} (s_0 \otimes I \otimes \beta_{\xi} \otimes I),$$

which contains the transition probabilities from $\Delta^{(0, \xi_1, \dots, \xi_m)}$ to $\Delta^{(\xi_1, \dots, \xi_m, \xi)}$, where $(\xi_1, \dots, \xi_m, \xi)$ means we rearrange elements in the vector in ascending order.

- If the last event is the arrival of a customer, we have

$$B_{00}^{(\xi_1, \dots, \xi_m)} = -(\beta_0 \otimes I) \left(S_0 \oplus \bigoplus_{j=1}^m S_{\xi_j} \right)^{-1} (s_0 \otimes \beta_{\xi} \otimes I), \quad (2.3)$$

which contains the transition probabilities from $\Delta^{(\xi_1, \dots, \xi_m)}$ to $\Delta^{(\xi_1, \dots, \xi_m, \xi)}$.

Example. For the PH/PH/2 queue, we have $\Omega_2(0) = \Delta^{(0)}$, $\Omega_2(1) = \Delta^{(1)} \cup \Delta^{(0,1)} \cup \Delta^{(0,2)}$, and $\Omega_2(2) = \Omega_2 = \Delta_0 \cup \Delta_1 \cup \Delta_2$. Note that we do not have the set $\Delta^{(2)}$ since an arriving customer selects the available server with the smallest index for service. The transition block matrices are given as

$$A_{0,1} = [B_{*0}^{(0)}, \mathbf{0}, \mathbf{0}]; \quad A_{1,2} = \begin{bmatrix} B_{00}^{(1)} & \mathbf{0} & \mathbf{0} \\ B_{*0}^{(0,1)} & \mathbf{0} & \mathbf{0} \\ B_{*0}^{(0,2)} & \mathbf{0} & \mathbf{0} \end{bmatrix};$$

$$A_{1,0} = \begin{bmatrix} B_{01}^{(1)} \\ B_{*1}^{(1)} \\ B_{*2}^{(2)} \end{bmatrix}; \quad A_{2,1} = \begin{matrix} \Delta_0 & \Delta^{(1)} & \Delta^{(0,1)} & \Delta^{(0,2)} \\ \Delta_1 & \mathbf{0} & B_{02}^{(1,2)} & B_{01}^{(1,2)} \\ \Delta_2 & \mathbf{0} & B_{22} & B_{21} \end{matrix}.$$

Transition blocks A_0 and A_2 are given in Eq. (2.2).

2.3. Queueing analysis

Analysis of the above Markov chain can lead to results on the queue length at arrival (or departure) time points. Let $\pi = (\pi_0, \pi_1, \pi_2, \dots)$ be the limiting probabilities of P . Let $\mathbf{p}_a = (p_{a,0}, p_{a,1}, p_{a,2}, \dots)$ be the limiting probabilities of the queue length just before the arrivals of customers. Then we have

$$p_{a,n} = \frac{1}{\lambda_0} \pi_n \begin{bmatrix} B_{00} \\ B_{10} \\ \vdots \\ B_{k0} \end{bmatrix}, \quad \text{for } n = k, k+1, k+2, \dots,$$

where λ_0 is the arrival rate of customers. Similarly, we can find the limiting probabilities of the queue length right after the departures of customers.

Note that the limiting probabilities can be computed using the matrix-analytic methods. For details see [2,7,8].

In order to obtain the waiting time distribution we can employ the well known absorbing Markov chain approach presented in several matrix analytic methods literature. The absorbing Markov chain can be created from the transition matrix equation (2.1). The limiting probabilities are then used to initiate the absorbing Markov chain.

3. The discrete time system model

3.1. The discrete time PH/PH/k

Let us now consider the case of the discrete PH/PH/k. This is more involved than its continuous time counterpart. Let us start

by considering the number of possible event occurrences. With k servers the number of event occurrences, N_k can be written as

$$N_k = \sum_{j=1}^{k+1} \binom{k+1}{j}.$$

This can be written as

$$N_k = N_{k-1} + 2^k, \quad k \geq 1, \quad \text{with } N_0 = 1.$$

The number of block matrices that have to be determined is $M_k = (N_k)^2$ for a k server problem. For example for the case of $k = 2$ we have $N_2 = 7$ and $M_2 = 49$. The seven points of events are given as

1. at arrivals: (a)
2. at service completions by server 1: (s_1)
3. at service completions by server 2: (s_2)
4. at service completions by server 1 and server 2: (s_1, s_2)
5. at points of joint arrival and service completion by server 1: (a, s_1)
6. at points of joint arrival and service completion by server 2: (a, s_2)
7. at points of joint arrival and service completion by server 1 and by server 2: (a, s_1, s_2).

The resulting Markov chain for the PH/PH/k is a GI/M/1 type with k sub-diagonals. Now we consider the block matrices of this GI/M/1 type Markov chain. Once again, let the PH-representations of the k services be (β_i, S_i) , $i = 1, 2, \dots, k$. Further let (β_0, S_0) be the PH-representation of the interarrival times. Let

$$\Omega = \{\sigma = (\sigma_0, \sigma_1, \dots, \sigma_k) : \sigma_i = 0 \text{ or } 1, i = 0, 1, 2, \dots, k\}.$$

Define $\mathbf{c}_\sigma = (c(0, \sigma_0), c(1, \sigma_1), \dots, c(k, \sigma_k))$, where

$$c(i, \sigma_i) = I, \text{ if } \sigma_i = 0; \quad \beta_i, \text{ if } \sigma_i = 1.$$

Define $\mathbf{d}_\sigma = (d(0, \sigma_0), d(1, \sigma_1), \dots, d(k, \sigma_k))$, where

$$d(i, \sigma_i) = S_i, \text{ if } \sigma_i = 0; \quad \mathbf{s}_i^0, \text{ if } \sigma_i = 1.$$

Define $A = (A_{\sigma^{(1)}, \sigma^{(2)}})$, where $\sigma^{(1)}, \sigma^{(2)} \in \Omega$, and

$$A_{\sigma^{(1)}, \sigma^{(2)}} = \mathbf{c}_{\sigma^{(1)}}(I - S_0 \otimes S_1 \otimes \dots \otimes S_k)^{-1} \mathbf{d}_{\sigma^{(2)}}.$$

Decompose A as $A = A_0 + A_1 + A_2 + \dots + A_k$ such that A_i contains all the blocks $A_{\sigma^{(1)}, \sigma^{(2)}}$ satisfying $\sum_{j=1}^k \sigma_j^{(2)} - \sigma_0^{(2)} + 1 = i$, for $i = 0, 1, 2, \dots, k$. In this way, we find all the $\{A_0, A_1, A_2, \dots, A_k\}$. Transition block matrices for boundary states can be found in a similar way.

Remark. Consider the conventional method that keeps track of the phases, i.e. ($I_a(t), I_1(t), I_2(t), \dots, I_k(t)$), where $I_a(t)$ refers to arrival phase at time t , and $I_j(t)$, $j = 1, 2, \dots, k$ refer to the phase of service of server j at time t . If we use that approach the number of phases is $\prod_{j=0}^k m_j$. The number of states of the new approach that we are proposing is $\prod_{j=0}^k (m_j + 1) - \prod_{j=0}^k m_j$. Thus, the new approach has a smaller state space if $\prod_{j=0}^k (1 + 1/m_j) < 2$.

3.2. Feasibility of this approach for the discrete time case

For the discrete PH/PH/k the advantage of studying the system at points of events may be questionable. The growth in N_k is exponential, so is the growth in M_k . We have $N_k = N_{k-1} + 2^k$, $k \geq 1$, hence generating the transition matrix in this case could be very involved. Once all the M_k matrices are determined generating their entries $A_{i_0, i_1, \dots, i_k; j_0, j_1, \dots, j_k}$ is straightforward since we know that

$$A_{i_0, i_1, \dots, i_k; j_0, j_1, \dots, j_k} = c(I - S_0 \otimes S_1 \otimes S_2 \otimes \dots \otimes S_k)^{-1} d,$$

and the values of c and d can be determined easily.

The reduction in matrix sizes only goes from $O(n^{k+1})$ to $O(n^k)$ when $m_j = n$, $\forall j$, if we use this approach.

4. Conclusion

Whereas there is a considerable saving in computational effort by studying the PH/PH/ k system with heterogeneous servers in continuous time according to the points of events, there may not be much savings, if any at all, for its discrete time counterparts. This is because the growth in the number of events occurring increases considerably as k increases as well as the number of matrices, while the reduction of the block matrix sizes goes from $O(n^{k+1})$ to $O(n^k)$. We might be better off using the traditional approach. And if the servers are identical then, of course, one may use the CSFP approach, presented in [4], to reduce the state space in both cases.

Acknowledgments

The authors acknowledge funding from the NSERC (Canada) and the NRF (South Africa) SARChI Chair programme cohosted by UP and CSIR.

References

- [1] A.S. Alfa, *Applied Discrete-Time Queues*, second ed., Springer, 2015.
- [2] Q.-M. He, *Fundamentals of Matrix-Analytic Methods*, Springer, New York, 2013.
- [3] Q.-M. He, A.S. Alfa, Construction of Markov chains for discrete time MAP/PH/K queues, *Perform. Eval.* 93 (2015) 17–26.
- [4] Q.-M. He, A.S. Alfa, Space reduction for a class of multidimensional Markov chains: A summary and some applications, *INFORMS J. Comput.* (2017) in press.
- [5] Q.-M. He, H. Zhang, Q.Q. Ye, An M/PH/K queue with constant impatient time, *Math. Methods Oper. Res.* (2017) accepted.
- [6] G. Latouche, V. Ramaswami, The PH/PH/1 queue at epochs of queue size changes, *Queueing Syst.* 25 (1997) 97–114.
- [7] G. Latouche, V. Ramaswami, Introduction to Matrix Analytic Methods in Stochastic Modelling, in: *SIAM-ASA Series on Statistics and Applied Probability*, SIAM, 1999.
- [8] M.F. Neuts, *Matrix-Geometric Solutions in Stochastic Models: An Algorithmic Approach*, The Johns Hopkins University Press, Baltimore, MD, 1981.
- [9] V. Ramaswami, Independent Markov proces in parallel, *Stoch. Models* 1 (1985) 419–430.
- [10] V. Ramaswami, D.M. Lucantoni, Algorithms for the multi-server queue with phase type service, *Stoch. Models* 1 (1985) 393–417.