

Variation of Access Control Policy for Two-queue Model

Cheng-Yuan Ku

Department of Information Management

National Chung Cheng University, Chia-Yi, Taiwan, ROC

E-mail: cooperku@mis.ccu.edu.tw or coopercy@ms16.hinet.net

Abstract

We consider access control problem for two multi-server loss queues in tandem, which is usually used to model networking systems. In previous work [4], it can be shown that under appropriate conditions the optimal managerial policy that maximizes the expected total discounted reward over an infinite horizon is given by a switching curve in the two-dimensional state space. In this paper, we discuss the variation of these switching curves with number of customers in system. Some experimental examples of unusual variation are presented. However, two sufficient conditions, which prevent this counterintuitive situation from happening were proposed.

Keywords: Variation of control policy, Access control, Tandem queues, Dynamic Programming, Optimal control

I. Introduction

Loss queueing models are often used to describe networking systems, especially for real-time applications and have been successful in analysis of resource allocation in these systems. Therefore, the optimal control problem for two-loss-queue-in-tandem model has attracted attention for years. Using the approach of dynamic programming method, the threshold-type managerial policy which reserves servers in second queue for internal customers coming from first queue has been extensively derived in [2,4]. From intuitive point of view, the optimal admission policy won't reserve more than one server for each customer at first station. However, some counterintuitive examples were found from experimental results. In these experiments, the system manager has to leave two servers open for one additional customer in the upstream queue according to the optimal managerial policy. A good reference on this topic can be found in [1,6]. The cause, which induces this counterintuitive situation, is still under investigation. However, two sufficient conditions, which prevent it from happening are derived according to the structure of optimality equation.

The two-queue model is presented in section 2. In section 3, we show two theorems, which describe the major topic in this paper and then experimental results are presented.

II. Model and Problem Formulation

Consider the system of multi-server loss queues in tandem as pictured in the following Fig. 1. The first queue A is a $M/M/m/m$ station in which λ_1 is the Poisson arrival rate, μ_A is the exponential service rate, R_1^A is the revenue for service to a new customer and P is the probability that customers leave system after service. The second queue B is a $M/n/n$ station in which λ_2 is the Poisson arrival rate, μ_B is the exponential service rate, R_2^B is the revenue for service to a new customer and R_1^B is the revenue for service to an internal customer from station A . Assume $R_1^B > R_2^B$, which means internal customers are more valuable than new customers at station B . Moreover, we also assume that $1 > P \geq 0$ and revenue is collected when a customer enters service. Suppose that only call admission control (CAC) is implemented to manage this system since it is of particular interest to high-speed networks. We consider the objective of maximizing the total discounted rewards over an infinite horizon. From the results in [4], it can be shown that the optimal admission policy doesn't control new customers at station A and internal customers at station B at all. Hence, the system manager only controls new customers at B .

This system can be modeled as a two-dimensional continuous-time Markov chain, with state $(i, j) \in \Omega : \{0, 1, \dots, m\} \times \{0, 1, \dots, n\}$ defined as the number of customers at stations A and B respectively. Uniformization results in an approximate discrete-time Markov chain by allowing fictitious transitions from a state to itself [5].

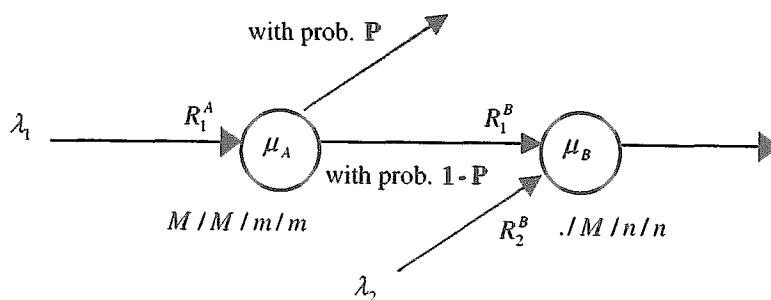


Fig. 1. Tandem multi-server loss queues

Choose an appropriate $Q \in \mathbb{R}$ s.t. $Q > \lambda_1 + \lambda_2 + m\mu_A + n\mu_B$ and let $p_1 = \frac{\lambda_1}{Q}$,

$$p_2 = \frac{\lambda_2}{Q}, \quad q_A = \frac{\mu_A}{Q}, \quad \text{and} \quad q_B = \frac{\mu_B}{Q}.$$

This discrete-time system has corresponding parameters p_1 , p_2 , q_A and q_B and the appropriate discount factor $\alpha < 1$.

Using the approach of dynamic programming, the following optimality equation, which need to be revised at the boundary states consists of all possible transitions with corresponding transition probabilities.

$$V(i, j) = \max_a \left\{ \begin{aligned} & p_1 [V(i+1, j) + R_1^A] \\ & + ap_2 [V(i, j+1) + R_2^B] \\ & + (1-P)(iq_A) [V(i-1, j+1) + R_1^B] \\ & + P(iq_A)V(i-1, j) + (jq_B)V(i, j-1) \\ & + (1-p_1 - ap_2 - iq_A - jq_B)V(i, j) \end{aligned} \right.$$

where $V(i, j)$ is the optimal value function and the optimal access control policy is given by $a: \Omega \rightarrow \{0,1\}$. $a(i, j) = 1$ (or 0) iff admit (or reject) the new customer at station B when the system is in (i, j) . The optimal admission policy admits a customer if the immediate revenue generated by that customer exceeds the expected loss in future discounted revenue caused by future blocking due to this customer. Thus we define the optimal difference function $\Delta(i, j) = V(i, j) - V(i, j+1)$. Then the optimal policy, in state (i, j) , admits a new customer at B , i.e. $a(i, j) = 1$; iff $\Delta(i, j) \leq R_2^B$.

III. Variation of Optimal Switching Curve

From results in [4], it can be shown that the optimal admission policy, which maximizes the expected total discounted reward over an infinite horizon is given by a switching curve in two-dimensional state space. The

system manager admits new customers at station B if state (i, j) is below the switching curve and rejects new customers while state (i, j) is on or above the threshold.

The managerial procedure should reserve servers for internal customers while the system tends to congestion according to the switching curve. In this paper, we focus on the relationship between the number of customers at station A and the number of reserved servers at station B for internal customers. Intuitively, the optimal policy won't reserve more than one server for each customer at station A . If it is correct, then the optimal switching threshold is not only nonincreasing in the number of customers at A but also is always horizontal or 45 degrees down as the switching curve 1 shown in Fig. 2.

However, a special experimental example shows the counterintuitive results. Adopting the following parameters for this system: $m = 6$, $n = 5$, $\alpha = 0.9999$, $Q = 100000$, $\lambda_1 = 18$, $\lambda_2 = 5$, $\mu_A = 18$, $\mu_B = 8$, $R_1^A = 100$, $R_1^B = 320$, $R_2^B = 10.5$ and $P = 0$, the optimal switching curve is pictured as the switching curve 2 in Fig. 2. From this curve, two interesting situations could be observed. First, the system manager reserves one space even though there is no customer at station A . That's because the arrival and service rates at station A are large, and the service rate at station B is relatively small. Therefore, the probability for new customers going through station A and then requiring service at station B in the near future is significantly large. In addition, the relative value of two types of customers at station B

$\frac{R_2^B}{R_1^B}$ is quite small. So the system manager values

internal customers much more than new customers at station B . Hence, one space is kept for this type of customers even while nobody is at A . Second, the system manager reserves three servers while there are two customers at A . However, he allocates five servers for upcoming internal arrivals while there are three customers at A . That means two more servers are reserved for one additional customer at the first station while the system is in this specific state. This phenomenon totally disagrees with our intuition.

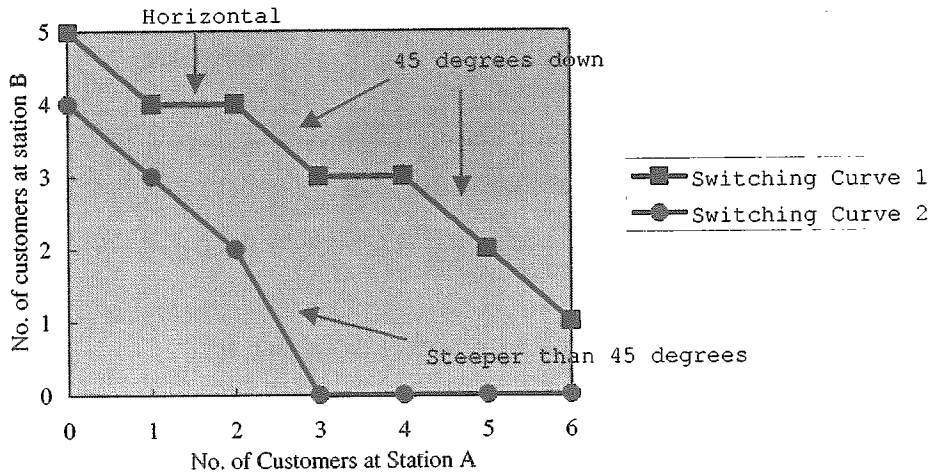


Fig. 2. The Optimal Switching Curve

From the structure of optimal difference function, it is easy to observe that if $\Delta(i, j) > \Delta(i+1, j-1)$ is true for any possible (i, j) , then the optimal switching curve is either horizontal or 45 degree down. In theorem 1, we can show that $\mathbb{P}q_A \geq q_B$ is a sufficient condition for $\Delta(i, j) > \Delta(i+1, j-1)$ for any possible (i, j) and therefore also a sufficient condition for the optimal policy to be either horizontal or 45 degrees down.

Theorem 1

- $\mathbb{P}q_A \geq q_B$
- $\Delta(i, j) > \Delta(i+1, j-1)$ for any possible (i, j)
- Optimal switching curve is either horizontal or 45 degrees down

In opposite, it indicates that if the leaving rate from station A is smaller and the flow rate from A to B is larger, then more servers should be reserved. Therefore, switching curve steeper than 45 degree is more possible. For the switching curve 2 in Fig. 2, the system with $\mathbb{P} = 0$, $\mu_A = 18$ and $\mu_B = 8$ has one segment of policy, which is steeper than 45 degrees. But if we let $\mathbb{P}q_A \geq q_B$ by setting $\mathbb{P} = 0.45 > \frac{8}{18} = 0.\bar{4}$, then the experimental result is pictured as one in Fig. 3.

$\mathbb{P} = 0.45$ satisfies $\mathbb{P}q_A \geq q_B$, and thus the switching curve is either horizontal or 45 degrees down. However, $\mathbb{P} = 0.3$ do not dissatisfy $\mathbb{P}q_A \geq q_B$, and yet the switching curve is still horizontal or 45 degrees down. Therefore, the condition $\mathbb{P}q_A \geq q_B$ is thus sufficient but not necessary.

The following theorem 2 shows that an alternative sufficient condition is $\frac{R_1^B}{R_2^B} \leq 1 + \frac{q_B}{p_1}$. However, it is of

note that this condition does not assure the diagonal monotonicity, $\Delta(i, j) > \Delta(i+1, j-1)$.

Theorem 2

$$\frac{R_1^B}{R_2^B} \leq 1 + \frac{q_B}{p_1} \quad \text{Optimal switching curve is either horizontal or 45 degrees down}$$

This sufficient condition is equivalent to $\frac{R_1^B}{R_2^B} \leq 1 + \frac{q_B}{p_1} \Leftrightarrow \frac{R_1^B - R_2^B}{R_2^B} \leq \frac{q_B}{p_1}$. Therefore, if $q_B = \frac{\mu_B}{Q}$ is going up, $p_1 = \frac{\lambda_1}{Q}$ is going down and

R_2^B is close to R_1^B , then less servers need to be reserved for the output flow from the first station. So, the optimal managerial policy is less possible to be steeper than 45 degree.

IV. Conclusion

From the experimental results, we can observe that the optimal access control policies for two-queue models are generally horizontal or 45 degrees down in the two-dimensional state space. A switching curve which is steeper than 45 degrees occurs rarely. This rare situation means the optimal managerial policy has to reserve more than one server at the second station for one additional customer at the first station in some specific states. Obviously, this result is counterintuitive. This phenomenon might be due to the procedure of uniformization but it is still under investigation. However, two sufficient conditions are presented to prevent this situation from happening. As for the omitted proofs can be found in [3].

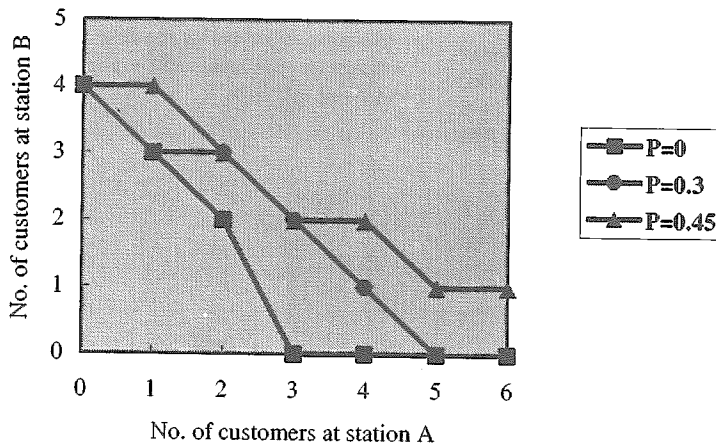


Fig. 3. The Optimal Policies for Adjusted P

References

- [1] E. Gelenbe, X. Mang and R. Onvural, "Bandwidth allocation and call admission control in high-speed network", *IEEE Communications Magazine*, vol. 35, pp. 122-129, May, 1997.
- [2] H. Ghoneim and S. Stidham, "Optimal control of arrivals to two queues in series", *European Journal of Operation Research*, 1985.
- [3] C.-Y. Ku, Access Control for Loss Network, Ph.D. dissertation, Department of EECS, Northwestern University, 1995.
- [4] C.-Y. Ku and S. Jordan, "Access control to two multiserver loss queues in series", *IEEE Trans. on Automatic Control*, vol. 42, pp. 1017-1023, July, 1997.
- [5] P. R. Kumar and P. P. Varaiya, *Stochastic Systems*, Prentice-Hall, 1986.
- [6] S. Stidham and R. Weber, "A survey of Markov decision models for control of networks of queues", *Queueing Systems*, vol. 13, pp. 291-314, 1993.