# State Dependent Social Optimization Where Customers Can Sell Their Place in Line

Avi Giloni

Sy Syms School of Business, Yeshiva University, BH-428, 500 West 185th St., New York, NY 10033, agiloni@yu.edu

Phil Troy

Les Entreprises TROYWARE, 5755 Eldridge, Cote St. Luc, Quebec, H4W 2E3, Canada, PhilTroy@Mail.com

## 1.   Introduction

Starting with Naor (1969), there has been a considerable amount of research on managing waiting in societal or organizational service facilities modeled as queues. This research includes that on state dependent control mechanisms, i.e. control mechanisms that depend on the number of customers waiting in the facility; it also includes research on priority based control mechanisms, i.e., control mechanisms that process jobs of customers with higher waiting costs first. To the best of our knowledge, there is no research that combines both approaches in a practical manner.

   Intuitively, it would seem that when internal customers incur different waiting costs, the "best" control mechanism would be a state dependent control combined with an ability to process customers with higher waiting costs first. This is because such a mechanism would make it possible to more tightly control the facility, and also make it possible to reduce waiting costs of those customers who incur them at the greatest rate. Unfortunately, while easy to specify mathematically, state dependent control mechanisms that group customers into different priority queues to achieve this level of control face at least two issues. The first issue is that the complexity of computing the parameters of these mechanisms grows extremely quickly in the number of possible jobs in each priority class and in the number of priority classes. The second and perhaps more important issue is that customers considering submitting their jobs to all but the highest priority class of service

face extremely high variances in the amount of time it will take for their job to get processed. This makes it hard for these customers to decide whether they should submit their job for processing. This in turn can result in facilities that process and/or admit almost entirely high priority class customers even when doing so is suboptimal, or in customers with lower waiting costs waiting excessively long times.

To overcome these issues, we extend Giloni and Troy's (2005) model for providing socially optimal state dependent control of service facilities modeled as queues with first come first served service, to permit customers with high waiting costs to swap positions for a negotiated fee with customers having lower waiting costs. This in turn makes it possible for customers having high benefit jobs and high waiting costs who arrive when the facility is busy to submit their jobs for processing, provided they can find a customer already in the queue with lower waiting costs who is willing to be paid to swap positions in the queue. Our approach also makes it possible for customers to control to their benefit the additional time they will have to wait. Since there is only a single queue, our approach has the additional benefit that computationally the problem reverts to that of a state dependent control with a single queue. We note that although there are many ways in which the negotiation process can take place, we only assume that the process will result in some customers with lower waiting costs being willing for a negotiated fee to swap places with a newly arrived customer with higher waiting costs.

Our research is strongly motivated by the literature on state dependent social optimization, initiated by Naor (1969) while examining societal service facilities modeled as queues, where customers waited for service. In that work, Naor shows that it can be socially suboptimal to allow customers to decide whether to submit requests for processing when doing so can cause additional waiting for subsequent customers. Mendelson (1985), in extending this analysis to facilities in which internal customers of a service facility only see average queue lengths, shows that the optimal toll reflects the opportunity costs that arriving customers cause for subsequent customers. Other related work includes that of Yechiali (1971), Knudsen (1972), and Lippman and Stidham (1977). Miller (1969) considers multiple server facilities that can be modeled as M/M/S/S queues where internal cus-

tomers having heterogeneous benefits are not allowed to wait. Whang (1986) considers single server facilities that can be modeled as M/M/1 queues where internal customers are allowed to wait, have the same linear waiting cost function, have processing requirements that are exponentially distributed with a common mean, and receive uniformly distributed benefits for processing of their jobs. Both Miller and Whang show that optimal control of the facilities they analyze can be achieved by decreasing job acceptance rates as the amount of work in the facility increases, and Whang achieves that control by using tolls. Giloni and Troy (2005) extend Whang's analysis by allowing for arbitrary waiting cost functions and job benefit distributions.

Our research is also closely related to state independent work by Mendelson and Whang (1990) and Affeche and Mendelson (2004). Mendelson and Whang develop a mechanism for setting prices that induce customers to join the appropriate queues controlled by state independent mechanisms. Affeche and Mendelson develop a mechanism whereby customers are induced to bid for their relative position in a single queue; they also show that while the resulting tolls are lower than the tolls computed without this bidding, the social benefits are higher.

In Section 2, we summarize the results of Giloni and Troy's (2005) state dependent control mechanism for a service facility that processes customers on a first come first served service facility, in context of an M/M/S/I queuing model. In section 3 we first show that extending this model to multiple priority classes very significantly increases the computational effort required to solve the problem and makes it very hard for customers to decide whether to submit jobs in all but the highest priority class. We then specify the new approach in which customers can sell their place in line. We show that the resulting problem makes it possible to achieve state dependent control with higher priority service for customers with higher waiting costs, while precluding the problems discussed above.

## 2. First Come First Served Service

Giloni and Troy (2005) propose a state dependent control mechanism for facilities in which internal customers are served on a first come first served basis. They show that the optimal control policy

is to set tolls equal to the opportunity costs inherent to the control problem, and that these opportunity costs, and thus the tolls, are non-decreasing in state.

In their model, $K$ groups of customers bring jobs to a service facility for processing. Group $k$ customers arrive at the facility at an average rate of $\lambda_k$ per unit time with exponentially distributed interarrival times. If they submit their job and it is accepted for processing, group $k$ customers receive a benefit $b_k$ after job processing is completed. While waiting for processing to be completed, group $k$ customers incur an expected waiting cost of $w_{i,k}$ that is non-decreasing and typically increasing in $i$, the number of jobs already in the facility when they arrive. They define $\beta_{i,k} = b_k - w_{i,k}$ to be the net benefit that group $k$ customers receive if their jobs are admitted to the facility when the facility is already in state $i$. Customers only submit jobs for processing if their net benefit is greater than or equal to any toll they are charged. The processing requirements of all jobs are exponentially distributed with a common mean. The facility is modeled as an $M/M/S/I$ queue having $S$ servers and a buffer that can hold $I - S$ jobs that processes all accepted jobs on a first come first served basis at the rate of $\mu_i$ jobs per unit time where $\mu$ is the average number of jobs processed by each server per unit time and $\mu_i = i\mu$ if $i \leq S$ and $S\,\mu$ otherwise.

Giloni and Troy observe that this problem is a Markov decision process, and that it can be formulated and solved as a non-discounted continuous time policy iteration problem. To do so, they specify the value determination equations for this problem (Howard 1960) which reduce to

$$\gamma = \sum_k \xi_{i,k} \lambda_k \left( \beta_{i,k} - \nabla v_i \right) + \mu_i \nabla v_{i-1}, \tag{1}$$

where $\gamma$ is the expected reward generated per unit time, $\nabla v_i = v_i - v_{i+1}$, and $\xi_{i,k}$ is the fraction of group $k$ jobs accepted when the facility is in state $i$.

## 3.  Priority Service

To adapt the first come first served model to allow for multiple priority classes without preemption, it is necessary to convert the scalar state variable $i$ to a state vector $\mathbf{i}$, and to explicitly state the before and after states involved in transitions and opportunity costs. Having done so, one can

observe that the optimal policy when the facility is in a particular state is to only admit customers whose net benefits are greater than the opportunity costs incurred by the facility in making a transition from its current state to the new state. The issue is with regard to the the calculations. When there is only a single class of service the value determination equations are tri-diagonal in the opportunity costs. In contrast, when there are multiple classes of service the value determination equations are no longer tri-diagonal, and the number of states is now approximately equal to the maximum number of customers that can wait in each class raised to the number of classes. Thus, if there are 3 priority classes each allowing up to 100 customers, it would require the solution of a non tri-diagonal set of equations with roughly 1,000,000 variables.

In addition to the increase in computational complexity, customers who join any but the highest priority class are faced with a very significant increase in the variability of their waiting times because they now not only have to deal with the variability of the processing times of the individuals in front of them but also with the variability in the number of customers that will join a higher priority class after they join their priority class queue. As an alternative, we allow customers with higher waiting costs to pay other customers with lower waiting costs to swap positions with them in a single queue.

To implement this policy it is only necessary to compute the expected arrival rate and net benefit of customers in each customer class that submit jobs when the facility is in a particular state and is charging a particular toll. We prove that these rates will always be at least as large as with the first come first served policy since customers can always obtain at least as much benefit with our priority mechanism than without it. To determine the actual rates while determining the optimal policy with policy iteration, we define the following notation. Let $\lambda_{i,k}$ be the average arrival rate to the system by customer group $k$ when there are $i$ customers in the system. Let $\beta_{i,k}$ be the average net benefit that customers from customer group $k$ receive when they enter the system when there are currently $i$ customers in the system including negotiation fees. With this notation, we propose the following procedure:

- Step 1: Propose an initial set of tolls for making the transition between adjacent states.

- Step 2: Simulate the arrival rates $(\lambda_{i,k})$ to the system and net benefit $(\beta_{i,k})$ rates of customers in each class for each state with these tolls under the assumption that negotiations (under some assumed negotiation process) will take place between arriving customers and customers already in the facility to induce them to swap positions.

- Step 3: Use the simulation results to compute the value determination equations for $\nabla v_i$.

- Step 4: Apply the policy improvement step.

- Step 5: Repeat the above steps until the policy converges.

THEOREM 1. *The above procedure converges to the optimal set of tolls.*

Sketch of Proof: In order to prove that the procedure converges, since all other steps of the procedure are standard policy iteration, all that is required is to ensure that the state dependent rates estimated from the simulations converge over the different iterations. To do so, it is only necessary to ensure that the estimation process includes the results from simulations in previous iterations, which in turn will ensure that as the number of iterations increases, the estimates will converge to their correct values, and thus the procedure will converge to the optimal set of tolls.

THEOREM 2. *The rate of net benefit per unit generated under this new policy is at least as good (and generally better) than that of the first come first served policy.*

Sketch of Proof: Consider the optimal first come first served policy. Apply policy improvement with the new policy of allowing users to switch places in each state (as per the new policy) and for new users who would not have joined because they were not able to switch places. If this results in an increase in $\gamma$ for even one state, the new policy is better than the previous policy. Since the old policy is a "feasible" solution under the new framework, it will not perform worse than before. Provided that there is at least one swap in all of the states, the new policy is actually better. Otherwise the old policy is optimal. □

### References

[1] Afeche, P., and Mendelson, H. 2004. Pricing and Priority Auctions in Queueing Systems with a Generalized Delay Cost Structure. *Management Science* 50, 7, 869-882.

[2] Giloni, A., and Troy, P. 2005. An Optimal Control Policy For Facilities Where Customers Wait For Service. working paper, Sy Syms School of Business, Yeshiva University.

[3] Howard, R.. 1960. *Dynamic Programming and Markov Processes*. The MIT Press.

[4] Knudsen, N.C. 1972. Individual and Social Optimization in a multi-server Queue with a General Cost-Benefit Structure, *Econometrica*, 40,3, 515-528.

[5] Lippman, S.A., and Stidham S. 1977. Individual versus Social Optimization in Exponential Congestion Facilities. *Operations Research.* 25,2, 233-247.

[6] Mendelson, H., and Whang, S., 1990. Optimal Incentive-Compatible Pricing for the M/M/1 Queue. *Management Science.* 38, 5, 870-883.

[7] Mendelson, H. 1985. Pricing Computer Services: Queueing Effects. *Communications of the ACM.* 28,3 312-321.

[8] Miller, B.L. 1969. A Queueing Reward System With Several Customer Classes. *Management Science.* 16,3 234-245.

[9] Naor, P. 1969. The regulation of queue size by levying tolls. *Econmetrica.* 37,1 15-24.

[10] Whang, S. 1986 The Value Of System State Information In Information Systems, working paper.

[11] Yechiali, U. 1971. On Optimal Balking Rules and Toll Charges in the GI/M/1 Queueing Process, *Operations Research.* 19, 349-370.