

Production Control with Backlog-Dependent Demand

Stanley B. Gershwin

Department of Mechanical Engineering
Massachusetts Institute of Technology
Cambridge, Massachusetts 02139-4307
USA
gershwin@mit.edu

Bariş Tan

Graduate School of Business
Koç University
Rumeli Feneri Yolu, Sarıyer
Istanbul, Turkey
btan@ku.edu.tr

Michael H. Veatch

Department of Mathematics
Gordon College
Wenham, MA 01984
USA
mike.veatch@gordon.edu

July 2, 2007

Abstract

We study a manufacturing firm that builds a product to stock to meet a random demand. Production time is deterministic, so that if there is a backlog, customers are quoted a lead time that is proportional to the backlog. In order to represent the customers' response to waiting, we introduce a *defection function* — the fraction of customers who choose not to order as a function of the quoted lead time. Unlike models with backorder costs, the defection function is related to customer behavior. Using a continuous flow control model with linear holding cost and Markov modulated demand, we show that the optimal production policy has a hedging point form. The performance of the system under this policy is evaluated, allowing the optimal hedging point to be found.

1 Introduction

1.1 Overview

THIS paper describes a model for determining a production policy for a factory that has impatient customers. Customers either obtain their products immediately when they make an

order (if there is finished product available) or they are given an estimate of how long they will have to wait for it. Some fraction of customers will find the waiting time estimate unacceptable and will cancel their orders. That fraction is an increasing function of the estimated waiting time.

The factory can be operated at any production rate up to its capacity. It incurs a profit for each item it sells, and it forgoes that profit when a customer cancels his order. It also incurs a cost for holding inventory. Demand is random, so the factory holds inventory when it can in order to limit future lost sales. The problem is to find the optimal production policy as a function of the current level of demand and of the current inventory or backlog.

We combine ideas from two areas of operations research: the literature on queues with impatient customers, and the literature on control-theoretical models of real-time production scheduling. The former deals with impatient customers who face servers operating in a fixed manner; the latter deals with customers who are always willing to wait for their product, but where the factory is controlled to minimize costs due to inventory and backlog. The former treats customers as discrete entities; the latter represents the product as a continuous material. To form the simplest model that combines both impatient customers and production control, we have chosen to model both the stream of customers and the product with continuous variables.

As in Kimemia and Gershwin (1983) and Bielecki and Kumar (1988), the difference between cumulative production and cumulative demand is called the *surplus*, represented by x . When x is positive, it is *finished goods inventory*. When it is negative, $-x$ is *backlog*¹. There are costs associated with inventory (including the interest cost on the raw material, the floor space devoted to storage, etc.).

If there is a positive inventory of finished goods, the customers make their purchases without delay and leave. If there is a backlog, customers are told their expected lead time, which is proportional to backlog. For any given lead time estimate, a fraction of the customers are willing to wait to make their purchases, but all others depart. The greater the lead time, the more customers leave without making a purchase. We assume that once the customer places an order, she does not withdraw her order.

We assume a perfectly reliable factory in which lead time is deterministic and proportional to backlog. The demand is random: it is Markov modulated with two levels. Extending the model to include machine unreliability would be straightforward. We have chosen not to do this in order to keep the notation simple and to keep our focus on the factory's response to customer behavior. In addition, this model is appropriate when the demand varies over a much longer time scale than machine downtimes.

We seek the optimal production rate as a function of the current surplus and the current demand level. The objective is to maximize long-run average profit, where revenue is diminished by customers who are not willing to wait for their products and the cost rate is linear in the surplus. Unlike the production control problem of Bielecki and Kumar (1988), there is no cost for backlog. We introduce a *defection function* $B(\cdot)$ where $B(x)$ is the fraction of customers who will not complete their orders when the backlog is x . Then the instantaneous demand is reduced by a factor of $1 - B(x)$. The defection function $B(\cdot)$ describes the delay tolerance or impatience of the customer

¹This definition is standard in production control and make-to-stock queues. It differs from backlog in a regular queue, which refers to the amount of work in the system. See the comment on *backlog* in Section 1.2.

population.

We only consider the effect of lead time on present sales. Lead time might have other consequences which we do not consider, such as a customer who finds the lead time too great being less likely to attempt to make a purchase in the future, and the damage to a firm's reputation when it has frequent large lead times. We do not consider the effect of price on customer behavior. Price is certainly important since reducing the price for long lead time deliveries can persuade some customers not to defect. Finally, we do not consider the service or prices offered by competitors.

We show that the structure of the solution is the hedging point policy of Kimemia and Gershwin (1983) and Bielecki and Kumar (1988). The manufacturing facility produces goods at the maximum production rate until the hedging point is reached. When x is at the hedging level, the production rate is set to the demand rate and the surplus remains constant. We assume that in the high demand state demand exceeds capacity; consequently, the maximum production rate is always used (except possibly when the system is in a transient state) when the demand is high.

When $B(\cdot)$ is piecewise constant, we can determine analytically the density function of the surplus under a hedging point policy. Since we can approximate any $B(\cdot)$ with a piecewise constant function, we can therefore solve systems with essentially any $B(\cdot)$. We express average profit as a function of the hedging point. Finally, we find the optimal hedging point numerically and evaluate the system performance measures including throughput (actual production rate), fill rate, and average backlog. Further analysis is given in Veatch (2007), where it is shown that the impact of $B(\cdot)$ on system performance can be expressed through a single parameter.

1.2 Related Work

Understanding how customers respond to waiting has been subject to numerous studies in consumer behavior literature, e.g. Taylor (1994). It is not desirable to make a customer wait for a service or a product and such a delay may result in losing the customer or receiving negative evaluation of the service or the product. If we consider the experience of an internet user as an example, it is observed that when users experience long waits for a web site's home page to load, they either quit using the web or redirect to an alternative web page (Weinberg 2000). Similarly, it is reported that the waiting times to load a web page affects evaluation of web sites (Dellaert and Kahn 1999).

Customers who demand a product have to wait unless that item is already available as a finished good. In a retailing or make-to-stock production setting, the effect of a stockout depends on the customer's personal commitment to the out-of-stock item and also on the difficulty of making a choice from the other available items. That is, some customers will wait for the item they originally sought; others will not wait for it, either because they buy a substitute or they do not buy anything (Fitzsimons 2000). Analysis of three years of purchasing behavior for 20,000 customers of a mail-order catalog revealed that when an item is out of stock, the percentage of customers who cancel their orders increases with the anticipated delay before the item is expected to ship (Anderson, Fitzsimons, and Simester 2003). In a make-to-order manufacturing setting, the quoted lead time and the price are the two most important decision criteria used in the negotiation process (Raiffa 1982). At a given price level, longer lead times may cause losing potential orders.

Customers who demand a service have to wait unless the server is idle and the service can start

immediately. In call centers, it is observed that the fraction of customers who abandon while they are on hold increases with the waiting time (Mandelbaum and Zeltyn 2004). Based on extensive statistical analysis of data from a telephone call center, Brown, Gans, Mandelbaum, Sakov, Shen, Zeltyn, and Zhao (2003) describe a survival function that is the fraction of customers who remain on hold after waiting for a given length of time. This survival function is conceptually very similar to the defection function used here and the data clearly shows that it decreases monotonically from 1.0 as the length of time increases. In health services, the elasticity of demand with respect to waiting time is estimated empirically using the waiting list for elective surgery in the UK (Martin and Smith 1999).

Developing analytical models that incorporate customer response to waiting directly has been the objective of numerous studies in the literature. There is a considerable literature on queues with impatient customers. Impatience in queues includes reneging (some customers abandoning the queue after waiting some time) and balking (some customers not joining the queue if the server is not immediately available). The reader is referred to a text book on queueing theory (such as Gross and Harris 1985 or Hall 1991) for an introduction to this literature and basic models.

The model described in this paper differs from the previous literature because its surplus/backlog (analogous to the queue length) is represented by a continuous quantity rather than an integer. Reneging does not occur in our model. We study a behavior which is analogous to balking: some customers do not complete an order if they find an excessive backlog (queue of orders) ahead of them. The most important way that our approach differs from the queueing theory literature is that we model the stream of customers and the product as continuous variables to seek a policy by which the production is controlled to maximize profit taking into account this balking behavior.

Since the 1980s, there has been an increasing interest in devising optimal production control policies that manage production in uncertain environments. An optimal flow-rate control problem for a failure prone machine subject to a constant demand was introduced by Olsder and Suri (1980) and Kimemia and Gershwin (1983). The single-part-type, single-machine problem was analyzed in detail by Bielecki and Kumar (1988). These papers penalized backlogs with a cost. Hu (1995) extended this work to the case in which no backlog is allowed and Martinelli and Valigi (2004) extended it to the case of limited backlog. Most of these studies assume a constant demand. A few consider optimal production control problems with random demand, as we do, including Fleming, Sethi, and Soner (1987), Ghosh, Araposthathis, and Markus (1993), Perkins and Srikant (2001), and Tan (2002). All of these papers model production as a continuous process (as we do here), and the optimal control for all of these systems was the hedging point policy (as we find for our problem). The most important difference between these papers and this one is that we allow backlog, but we penalize backlog with a loss of demand to incorporate the customer response to waiting directly in our model.

Note that the term *backlog* has two different meanings in the two different theoretical areas that we draw on. In the queueing literature, backlog is the amount of work in the system; in the production control literature, it is the negative part of the difference between cumulative production and cumulative demand. We use it in the latter sense in this paper. From the point of view of a customer, whether a delay is caused by work in progress ahead in the queue (backlog in the queueing sense) or by the production system falling behind demand due to demand fluctuations (or

capacity variability) is indistinguishable. (In fact, the former is caused by the latter.) Thus, the two meanings of backlog are unified here.

1.3 Outline

The model and its assumptions are described in Section 2, where backlog-dependent demand is introduced and the production control problem is stated. The policy that maximizes the profit is characterized in Section 3. The model is analyzed and the steady state probability distributions are formulated in Section 4. Section 5 describes the evaluations of the objective function and of other performance measures of interest. To illustrate the behavior of the model, one family of defection functions is studied in Section 6. Results are sensitive to both the median delay tolerance and the variability of the delay tolerance. Section 7 contains a summary of the paper and several proposed research directions.

2 Model Description

2.1 Basic Model

We consider a make-to-stock system with a single manufacturing facility that produces to meet the demand for a single item. Production, demand, inventory, and backlog are all represented by continuous (real) variables. The demand rate at time t is denoted by $d(t)$. This demand rate has two possible values: high demand (μ_H) and low demand ($\mu_L, \mu_L < \mu_H$). It is convenient to define a demand state $D(t)$ such that $d(t) = \mu_H$ when $D(t) = H$, and $d(t) = \mu_L$ when $D(t) = L$.

The times to switch from the high demand state to the low demand state and from the low demand state to the high demand state are assumed to be exponentially distributed random variables with rates λ_{HL} and λ_{LH} . This model is suitable for describing demand which is stationary in the long run, but whose mean shifts temporarily as a result of promotions, competitor actions, etc. The time since the last state change does not change the expected time until the next state change.

The maximum production rate of the manufacturing facility is \underline{u} . The actual production rate of the manufacturing facility at time t is a control variable which is denoted by $u(t)$, $0 \leq u(t) \leq \underline{u}$. We assume that the production capacity \underline{u} is sufficient to meet the demand when it is low but insufficient when it is high, i.e., $\mu_L < \underline{u} < \mu_H$. (Note that if $\underline{u} > \mu_H$, the problem is trivial: it is always possible to keep x at 0 and therefore a backlog situation never happens. Similarly, if $\underline{u} < \mu_L$, the problem is also trivial: the manufacturing facility is run at the maximum rate all the time.)

Not all the demand results in orders because some customers are discouraged by backlog and they defect. We define $x(t)$ to be the production surplus, the difference between cumulative production and cumulative orders during $[0, t]$. If $x(t)$ is positive, it is the finished goods inventory; if it is negative, $-x(t)$ is the backlog.

The reward per unit for goods produced in the factory is A and the inventory carrying cost is g^+ (dollars per unit per time). As indicated earlier, we do not include the corresponding backlog cost g^- , which does appear in Bielecki and Kumar (1988) and many other papers.

2.2 Backlog-Dependent Demand

When there is backlog (i.e., when $x < 0$), we define $B(x, D)$ to be the fraction of potential customers who choose not to order when the backlog is x and the demand is D . The *defection function*, $B(\cdot, \cdot)$, satisfies

$$\left. \begin{aligned} 0 \leq B(x, D) \leq 1 \text{ for all } x, D, \\ x > 0 \implies B(x, D) = 0 \text{ for all } D, \\ B(0, L) = 0, \\ B(0, H) = \min \left\{ \frac{\mu_H - u}{\mu_H}, B(0^-, H) \right\}, \\ \text{For some } x, \text{ for all } D, B(x, D) \geq \frac{\mu_H - u}{\mu_H}. \end{aligned} \right\} \quad (1)$$

The first condition is required by the definition of B as a fraction. The second says that potential customers are never motivated to defect when there is finished goods inventory.

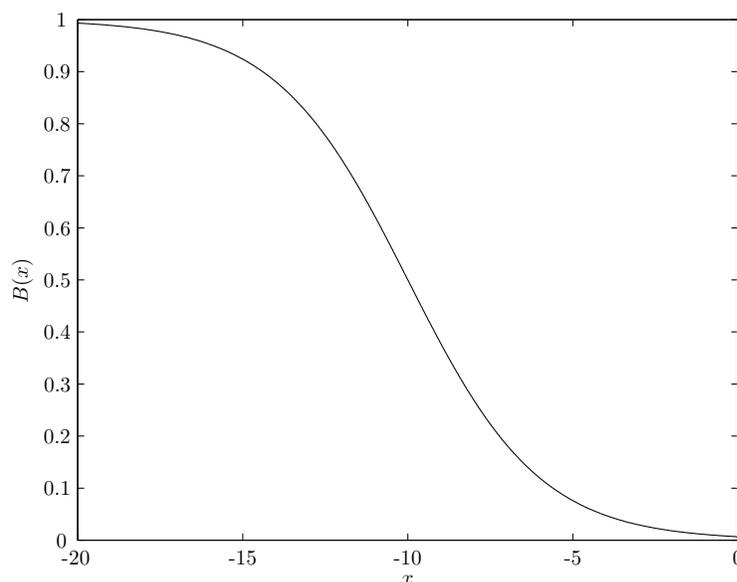
The third and fourth condition prevent excess defection when $x = 0$. Only customers with no patience and who cannot be served immediately should defect. Thus, when $D = L$ no customers defect and when $D = H$ the fraction defecting is the minimum of the fraction with no patience and the fraction that could not be served immediately. Note that we could have set $B(0) = 0$, avoiding the dependence of B on D , but this would make chattering possible in state $(0, H)$. Since the dependence on D is merely a technical convenience, in the remainder of the paper we suppress the dependence on D in our notation, and write $B(x)$.

The last condition guarantees that $x(t)$ is bounded from below. It is implied by the natural condition $\lim_{x \rightarrow -\infty} B(x) = 1$, which says that nobody is infinitely patient. This condition is further discussed at the end of Section 2.3.

If $B(x)$ is a non-increasing function of x , we say that B is a *monotonic defection function*. If B is monotonic, more customers are impatient if there is a longer wait. In the following, we restrict our attention to monotonic defection functions. Lastly, we assume that $B(x)$ is right continuous left limit for $x < 0$. Figure 1 shows an example of a $B(\cdot)$ function that satisfies these conditions.

In the Appendix of Tan and Gershwin (2001), the customer defection function that is generated by customers who choose the shortest of two queues is derived. (In this case, customers do not join queue i ($i = 1, 2$) if they see an opportunity to wait less at queue j ($j = 1, 2; j \neq i$); if the current length of queue i is n , the customer defection function² $B^*(n)$ is the probability that the length of queue j is less than n .) This preliminary analysis yields a $B(\cdot)$ function similar to the one depicted in Figure 1. More specifically, Tan and Gershwin (2001) consider two M/M/1 queues operating in parallel. An arriving customer chooses the queue that has the smallest number of customers (if the queue lengths are equal, she chooses one of them with equal probability). In

²In this case the customer defection function is a function of the queue length instead of a backlog.

Figure 1: A sample $B(\cdot)$ function

the long run, each server sees that an arriving customer may decide not to join her queue (i.e., to defect) with a probability depending on the number of waiting customers upon the arrival of this customer. This probability as a function of the queue length is the discrete version of the defection function considered here. It is shown that this function is a non-decreasing function of the number of customers in the system.

The empirical study of Anderson, Fitzsimons, and Simester (2003) also supports the kind of the defection behavior described here.

Note that $B(x) = 1$ for all $x < 0$ corresponds to the lost sales case. In this situation, no customers are willing to wait to receive their goods. All revenues are lost whenever there is any backlog.

When the surplus level is $x < 0$, the time until the next arriving customer order will be completed, i.e., the production lead time, is $-x/\underline{u}$. This is the time to clear the current backlog, assuming $u = \underline{u}$ until the backlog is cleared. The lead time dependent demand case can therefore be treated by using $B(x) = B(-\underline{u}t_q) = \tilde{B}(t_q)$ as the fraction of potential customers that choose not to order when the quoted lead time is $t_q = -x/\underline{u}$.

2.3 Production Control Problem

The decision variable is the rate at which the goods are produced at the plant at time t , $u(t)$. The dynamics of $x(t)$, on any sample path, are given by

$$\frac{dx(t)}{dt} = u(t) - d(t)(1 - B(x(t))).$$

at points where the derivative exists.

In earlier versions of the production control problem (for example Kimemia and Gershwin 1983; Bielecki and Kumar 1988), the $(1 - B(x(t)))$ factor did not appear. (That is, B was always 0.) It is present in this version to represent the lost demand due to customer balking. The quantity $d(t)(1 - B(x(t)))$ is the rate at which sales actually take place, and is therefore the rate at which the surplus/backlog is diminished by sales.

The revenue rate at time t is therefore $Ad(t)(1 - B(x(t)))$. The profit rate is the difference between the revenue generated by sales (i.e., realized demand) and the inventory carrying costs, which are assumed to be linear:

$$g(x(t), D(t)) = Ad(t)(1 - B(x(t))) - g^+x^+(t).$$

In a slight abuse of notation, we write $u = u(t)$ to denote the policy. The expected profit $J^u(x, D; T)$ of policy u in the interval $[0, T]$ from the initial state $x(0) = x, D(0) = D$, is defined as

$$J^u(x, D; T) = E_{x,D} \int_0^T g(x(t), D(t)) dt,$$

where $E_{x,D}$ denotes expectation with respect to the initial state $x(0) = x, D(0) = D$. Similarly, the long-run average profit of policy u is defined as

$$\bar{J}^u = \limsup_{T \rightarrow \infty} \frac{1}{T} J^u(x, D; T).$$

The production control problem is

$$\bar{J}^* = \max_u \bar{J}^u \quad (2)$$

subject to

$$\frac{dx(t)}{dt} = u(t) - d(t)(1 - B(x(t))) \quad (3)$$

$$0 \leq u(t) \leq \underline{u} \quad (4)$$

$$d(t) = \begin{cases} \mu_H & \text{if } D(t) = H \\ \mu_L & \text{if } D(t) = L \end{cases} \quad (5)$$

$$\left. \begin{aligned} \mathbf{prob}(D(t + \delta t) = H | D(t) = L) &= \lambda_{LH} \delta t + o(\delta t), \\ \mathbf{prob}(D(t + \delta t) = L | D(t) = H) &= \lambda_{HL} \delta t + o(\delta t). \end{aligned} \right\} \quad (6)$$

Define $N(t)$ to be the cumulative demand during $[0, t]$. The average demand rate is

$$Ed = \lim_{t \rightarrow \infty} \frac{E[N(t)]}{t} = \mu_H e + \mu_L(1 - e)$$

where $e = \lambda_{LH}/(\lambda_{HL} + \lambda_{LH})$ is the percentage of the time the demand is high. The system is stable under some policy if $Ed(1 - \lim_{x \rightarrow -\infty} B(x)) < \underline{u}$ or $\lim_{x \rightarrow -\infty} B(x) > \frac{(Ed - \underline{u})}{Ed}$. In fact, there is no reason for $u(t)$ to be anything but maximal when $x(t)$ is negative. Then (1) implies that $x(t)$ is bounded from below. To see this, let $x^* = \max\{x : \underline{u} - \mu_H(1 - B(x)) \geq 0\}$ and assume that $x(0) > x^*$. From (3), $dx(t)/dt \geq 0$ for all $x(t) \leq x^*$, for $D(t) = L$ or H (since $\mu_H > \mu_L$). Then $x(t) \geq x^*$ for all t . Even if Ed is greater than \underline{u} , enough impatient customers will defect to guarantee that $x(t)$ is bounded from below.

3 Characterization of the Policy

The solution of problem (2)–(6) is a stationary feedback control $u(t) = u(x(t), D(t))$ that satisfies the Bellman equation. Define the differential cost (actually reward) for policy u as

$$V^u(x, D) = \lim_{T \rightarrow \infty} J^u(x, D; T) - T\bar{J}^u$$

and the differential cost for the optimal policy by V when these limits exist. The differential cost V satisfies the *maximum principle*, which asserts that, assuming that $\partial V/\partial x$ exists,

$$\bar{J}^* = \max_u \left\{ -g(x, L) + \frac{\partial V}{\partial x}(x, L)(u - \mu_L(1 - B(x))) + [V(x, H) - V(x, L)]\lambda_{HL} \right\} \quad (7)$$

for $D = L$, and

$$\bar{J}^* = \max_u \left\{ -g(x, H) + \frac{\partial V}{\partial x}(x, H)(u - \mu_H(1 - B(x))) + [V(x, L) - V(x, H)]\lambda_{LH} \right\} \quad (8)$$

for $D = H$. The maximizations are taken over constraints (4).

The differentiability of V in intervals of constant control follows from the general theory of jump Markov processes in Rishel (1975). Differentiability where the control changes could be proven using the methods of Sethi, Suo, Taksar, and Zhang (1997). However, we can avoid this issue by replacing $\partial V/\partial x$ by the appropriate one-sided derivative, namely, $\partial^+ V/\partial x$ when $\frac{dx}{dt}(t^+) \geq 0$ and $\partial^- V/\partial x$ when $\frac{dx}{dt}(t^+) < 0$. The existence of one-sided derivatives follows from Rishel's theory.

The optimal policy has a very simple form.

Theorem 1 *The optimal policy has hedging point form: For some $Z(L) \geq 0$,*

$$u^*(x, D) = \begin{cases} \mu_L & \text{if } x = Z(L) \text{ and } D = L \\ \underline{u} & \text{if } x < Z(L). \end{cases}$$

A proof is given in the Appendix, based on comparing the problem to the lost sales unreliable machine problem. The policy on the transient states $x > Z(L)$ does not affect the average cost.

However, it seems apparent from similar problems that the optimal policy on these states uses thresholds $0 \leq Z(L) \leq Z(H)$, and that the complete policy is

$$u^*(x, D) = \begin{cases} 0 & \text{if } x > Z(L) \text{ and } D = L \\ \mu_L & \text{if } x = Z(L) \text{ and } D = L \\ \underline{u} & \text{if } x < Z(L) \text{ and } D = L \\ \\ 0 & \text{if } x > Z(H) \text{ and } D = H \\ \underline{u} & \text{if } x \leq Z(H) \text{ and } D = H. \end{cases}$$

Although $Z(L)$ is a hedging point, $Z(H)$ is not; this is because $dx(t)/dt < 0$ at $x(t) = Z(H)$, $D(t) = H$. That is, it is not possible for $x(t)$ to remain at $Z(H)$.

4 Analysis of the Model

In this section, we analyze the model with backlog-dependent demand. We calculate the steady-state probability distribution of x and D assuming that the system is operated under the policy of Section 3. In Section 5, we evaluate the expected profit (as well as other performance measures). Then we find the optimal policy by finding the value of $Z(L)$ that maximizes the expected profit.

4.1 Dynamics

The analysis of even simple systems with general non-zero $B(x)$ results in non-closed form solutions. In order to treat a wide variety of backlog-dependent demand functions, it is convenient to assume that $B(x)$ is piecewise constant. That is,

$$B(x) = \begin{cases} 0 & x > 0, \\ B_1 & 0 \geq x > \beta_1, \\ B_i & \beta_{i-1} \geq x > \beta_i \quad i = 2, \dots, M. \end{cases} \quad (9)$$

where $\beta_i, B_i, i = 1, \dots, M$ are constants, $0 > \beta_i > \beta_{i+1}$. From (1) and the monotonicity of $B(\cdot)$,

$$0 \leq B_i < B_{i+1} \leq 1.$$

By a proper choice of these constants, and for large enough M , any monotonic $B(x)$ can be arbitrarily closely approximated.

It is also convenient to include the lower bound on x in the discretization of $B(x)$. In order to analyze only the recurrent states, B_M can be chosen in such a way that it satisfies

$$\underline{u} = \mu_H(1 - B_M)$$

Let us define region boundaries $R_0 > R_1 > \dots > R_J$ with

- for $Z(L) > 0$: $R_0 = Z(L)$, $R_1 = 0$, $R_i = \beta_{i-1}$ for $i = 2, \dots, M + 1$ and $J = M + 1$,

- for $Z(L) = 0$: $R_0 = 0$, $R_i = \beta_i$ for $i = 1, \dots, M$, and $J = M$.

The recurrent values of x are $[R_J, R_0]$. Call (R_{i+1}, R_i) region i . Within each of these regions, the right side of the x dynamics is constant and $g(x, D)$ is linear.

Let Δ_i^L be the rate of change of x in region i when the demand state D is low (L). Then

$$\Delta_i^L = u - \mu_L(1 - B(x)), \quad R_i < x < R_{i+1}, \quad i = 0, 1, \dots, J - 1 \quad (10)$$

where u and $B(x)$ are constant in each region, so Δ_i^L is constant.

Similarly, let Δ_i^H be the rate of change of x in region i when $D = H$. Then

$$\Delta_i^H = u - \mu_H(1 - B(x)), \quad R_i < x < R_{i+1}, \quad i = 0, 1, \dots, J - 1 \quad (11)$$

where u is described in Section 3.

A sample path of a system where $B(x)$ is given as a five-level step-wise constant function is depicted in Figure 2. When the demand is high, $u = \underline{u}$. However, since $\underline{u} < \mu_H$, the surplus decreases with rate $\Delta_0^H = \underline{u} - \mu_H$ for $x > 0$ (Region 0). For $x < 0$, $B(x)$ increases step-by-step from β_1 to β_5 . When $x = \beta_5$, the demand rate of the customers who have not defected is equal to the maximum production rate and therefore, the backlog stays at the lower bound until the demand switches to low. When the demand is low, x increases with a step-wise constant slope until it reaches the hedging level at $Z(L)$.

In the following sections, we describe how the optimal policy is determined. First, the system is evaluated by determining the probability density functions in the interior, and probability masses at the upper and lower levels for a given value of $Z(L)$. Then, the optimal value of the hedging level is determined by maximizing the expected profit.

4.2 Probability distribution

When the surplus/backlog x is not equal to the upper or lower levels (R_0 or R_J), the system is said to be in the *interior*. The system state at time t is $(x(t), D(t))$ where $R_J < x(t) < R_0$ and $D(t) \in \{H, L\}$.

The time-dependent system state probability distribution for the interior region, $F_D(t, x)$, is defined as

$$F_D(t, x) = \mathbf{prob}[D(t) = D, x(t) \leq x], \quad t \geq 0, \quad D \in \{H, L\}, \quad R_J < x(t) < R_0. \quad (12)$$

The time-dependent system state density functions are defined as

$$f_D(t, x) = \frac{\partial F_D(t, x)}{\partial x} \quad t \geq 0, \quad D \in \{H, L\}, \quad R_J < x(t) < R_0. \quad (13)$$

It is possible to show that the process is ergodic by observing that in the Markov process model, all of the states constitute a single communicating class. It is also possible to demonstrate aperiodicity. Thus, steady-state density functions exist. They are

$$f_D(x) = \lim_{t \rightarrow \infty} f_D(t, x), \quad D \in \{H, L\}, \quad R_J < x(t) < R_0. \quad (14)$$

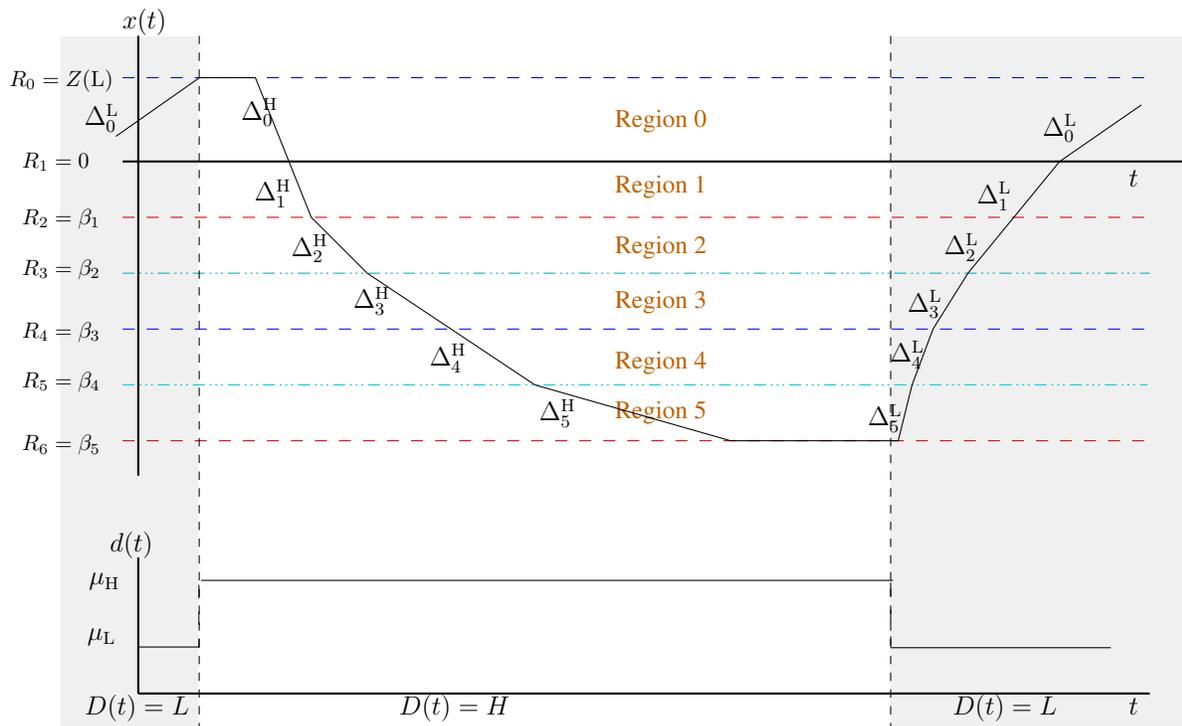


Figure 2: Sample path for a system

4.3 Region i : $R_{i+1} < x < R_i$

Suppose $R_{i+1} < x(t + \delta t) < R_i$, $i = 1, 2, \dots, J - 1$, and $D(t + \delta t) = H$. Then, since we are modeling this system as a Markov process,

$$f_H(t + \delta t, x) = f_H(t, x - \Delta_i^H \delta t)(1 - \lambda_{HL} \delta t) + f_L(t, x)(\lambda_{LH} \delta t) + o(\delta t) \quad (15)$$

where $o(\delta t)$ approaches to zero faster than δt . This equation can be written in differential form, for $\delta t \rightarrow 0$, as

$$\frac{\partial f_H(t, x)}{\partial t} + \Delta_i^H \frac{\partial f_H(t, x)}{\partial x} = -\lambda_{HL} f_H(t, x) + \lambda_{LH} f_L(t, x). \quad (16)$$

Taking the limit of (16) as $t \rightarrow \infty$ yields the following steady-state differential equation for $f_H(x)$:

$$\Delta_i^H \frac{df_H(x)}{dx} = -\lambda_{HL} f_H(x) + \lambda_{LH} f_L(x), \quad R_{i+1} < x < R_i. \quad (17)$$

Following the same steps for f_L yields

$$\Delta_i^L \frac{df_L(x)}{dx} = \lambda_{HL} f_H(x) - \lambda_{LH} f_L(x) \quad R_{i+1} < x < R_i. \quad (18)$$

In order to solve the set of first order differential equations given in (17) and (18), two boundary conditions for each region are needed. First, note that at any given level of the finished goods inventory, the number of upward crossings must be equal to the number of downward crossings. Let $\nu(D, \xi, T)$ denote the total number of level crossings in demand state D , at surplus level ξ , in the time interval $[t, t + T]$. Then

$$\lim_{T \rightarrow \infty} \nu(H, \xi, T) = \lim_{T \rightarrow \infty} \nu(L, \xi, T). \quad (19)$$

Renewal analysis shows that

$$\lim_{T \rightarrow \infty} \frac{\nu(D, \xi, T)}{T} = \Delta_i^D f_D(\xi) \quad (20)$$

where Δ_i^D is the rate of change in the buffer level when the demand state is D and ξ is in region i , and $f_D(\xi)$ is the steady-state density function. This kind of analysis was also employed by Yeralan and Tan (1997). Then, equation (19) can be written as

$$-\Delta_i^H f_H(x) = \Delta_i^L f_L(x). \quad (21)$$

Using this result in equation (17) gives the following first order differential equation

$$\frac{df_H(x)}{dx} = \left(-\frac{\lambda_{HL}}{\Delta_i^H} - \frac{\lambda_{LH}}{\Delta_i^L} \right) f_H(x) \quad (22)$$

whose solution is

$$f_H(x) = c_i e^{\eta_i x}, \quad R_{i+1} < x < R_i \quad (23)$$

where

$$\eta_i = -\frac{\lambda_{HL}}{\Delta_i^H} - \frac{\lambda_{LH}}{\Delta_i^L}$$

and c_i is a constant to be determined. Following equation (21),

$$f_L(x) = -c_i \frac{\Delta_i^H}{\Delta_i^L} e^{\eta_i x}, \quad R_{i+1} < x < R_i. \quad (24)$$

4.4 External Boundary Conditions

The steady-state probabilities P^0 and P^J that the finished goods inventory is equal to the hedging level $Z(L)$ and the lowest level \underline{X} are defined as

$$P^0 = \lim_{t \rightarrow \infty} \mathbf{prob}[x(t) = R_0], \quad (25)$$

$$P^J = \lim_{t \rightarrow \infty} \mathbf{prob}[x(t) = R_J]. \quad (26)$$

The inventory level can increase only when the demand is low. When x increases and reaches the level $R_0 = Z(L)$, the inventory level stays at the upper level until the demand rate increases to μ_H and x starts decreasing. The expected remaining time for the state of the demand to change from L to H is $1/\lambda_{LH}$. Then P^0 is fraction of time that $x = Z(L)$:

$$P^0 = \lim_{T \rightarrow \infty} \frac{\nu(L, R_0, T)}{T} \frac{1}{\lambda_{LH}} = \Delta_0^L f_L(R_0) \frac{1}{\lambda_{LH}} = -c_0 \frac{\Delta_0^H}{\lambda_{LH}} e^{\eta_0 R_0}. \quad (27)$$

Define

$$\nu(j, R_i^+, T) = \lim_{\substack{\xi \rightarrow R_i \\ \xi > R_i}} \nu(j, \xi, T); \quad \nu(j, R_i^-, T) = \lim_{\substack{\xi \rightarrow R_i \\ \xi < R_i}} \nu(j, \xi, T)$$

Then, similarly,

$$P^J = \lim_{T \rightarrow \infty} \frac{\nu(H, R_J^+, T)}{T} \frac{1}{\lambda_{HL}} = -\Delta_{J-1}^H f_H(R_J) \frac{1}{\lambda_{HL}} = -c_{J-1} \frac{\Delta_{J-1}^H}{\lambda_{HL}} e^{\eta_{J-1} R_J}. \quad (28)$$

Let us also define P_i^H and P_i^L $i = 0, 1, \dots, J-1$ as the probabilities that the process is in region i in the long run when the demand is high and when it is low, respectively:

$$P_i^H = \lim_{t \rightarrow \infty} \mathbf{prob}[R_i < x(t) < R_{i+1}, D(t) = H] \quad i = 0, \dots, J-1, \quad (29)$$

$$P_i^L = \lim_{t \rightarrow \infty} \mathbf{prob}[R_i < x(t) < R_{i+1}, D(t) = L] \quad i = 0, \dots, J-1. \quad (30)$$

Once the density functions are available, P_i^H and P_i^L can be evaluated as

$$P_i^H = \int_{R_{i+1}}^{R_i} f_H(x) dx \quad i = 0, \dots, J-1, \quad (31)$$

$$P_i^L = \int_{R_{i+1}}^{R_i} f_L(x) dx \quad i = 0, \dots, J-1. \quad (32)$$

4.5 Internal Boundary Conditions

To complete the derivation of the density functions, the coefficients c_i , $i = 0, 1, \dots, J-1$ must be determined. Since there are J unknowns, J boundary conditions are needed. The $J-1$ internal boundary conditions come from the equality of the number of upward and downward crossings at the region boundaries. For large T ,

$$\lim_{T \rightarrow \infty} \nu(j, R_i^+, T) = \lim_{T \rightarrow \infty} \nu(j, R_i^-, T), j \in \{H, L\}, i=1, 2, \dots, J-1. \quad (33)$$

By using equation (20), this equation can be written

$$\Delta_{i-1}^j f_j(R_i^+) = \Delta_i^j f_j(R_i^-), j \in \{H, L\}, i=1, 2, \dots, J-1. \quad (34)$$

5 Solution of the Model

5.1 Coefficients

Writing (34) in terms of the solution of the density function for $j = H$ given in equation (23) yields

$$\Delta_{i-1}^H c_{i-1} e^{\eta_{i-1} R_i} = \Delta_i^H c_i e^{\eta_i R_i}, i=1, 2, \dots, J-1, \quad (35)$$

or

$$c_i = \frac{\Delta_{i-1}^H}{\Delta_i^H} e^{(\eta_{i-1} - \eta_i) R_i} c_{i-1}, i = 1, 2, \dots, J-1 \quad (36)$$

Then all the constants c_i $i = 1, 2, \dots, J-1$ can be determined by c_0 , since $c_i = \phi_i c_0$ where

$$\phi_i = \prod_{j=1}^{i-1} \frac{\Delta_{j-1}^H}{\Delta_j^H} e^{(\eta_{j-1} - \eta_j) R_j}, i = 1, \dots, J-1 \quad (37)$$

and $\phi_0 = 1$.

Finally, the constant c_0 is determined by using the normalizing condition. The sum of all the probabilities must add up to 1, or

$$P^0 + \sum_{i=0}^{J-1} (P_i^H + P_i^L) + P^J = 1. \quad (38)$$

Equations (23), (24), (27), (28), (31) and (32) yield

$$c_0 = \left[\frac{(\mu_H - \underline{u})e^{\eta_0 R_0}}{\lambda_{LH}} + \sum_{i=0}^{J-1} \phi_i \frac{(\Delta_i^L - \Delta_i^H)}{\Delta_i^L} X_i - \frac{\phi_{J-1} \Delta_{J-1}^H e^{\eta_{J-1} R_J}}{\lambda_{HL}} \right]^{-1} \quad (39)$$

where

$$X_i = \begin{cases} (e^{\eta_i R_i} - e^{\eta_i R_{i+1}})/\eta_i & \text{if } \eta_i \neq 0, \\ (R_i - R_{i+1}) & \text{if } \eta_i = 0. \end{cases}$$

5.2 Evaluation of the Objective Function

In order to determine the optimal values of the hedging levels, the profit must be evaluated. Let Π be the average revenue rate generated by the plant. It is given by

$$\Pi = \lim_{T \rightarrow \infty} E \left[\frac{1}{T} \int_0^T Ad(\tau)(1 - B(x(\tau)))d\tau \right].$$

From (3), this is

$$\lim_{T \rightarrow \infty} E \left[\frac{1}{T} \int_0^T A \left(u - \frac{dx}{dt} \right) d\tau \right] = \lim_{T \rightarrow \infty} E \left[\frac{1}{T} \int_0^T Aud\tau - \frac{A}{T}(x(T) - x(0)) \right].$$

But we know that $E[x(T) - x(0)]$ will be bounded since the system is stable. Since the difference between the cumulative production and cumulative demand is finite in the long run, the profit term in the objective function can also be written using the production rate rather than the demand rate. Therefore,

$$\Pi = \lim_{T \rightarrow \infty} E \left[\frac{1}{T} \int_0^T Aud\tau \right] = A (\mathbf{prob}[x < R_0] \underline{u} + \mathbf{prob}[x = R_0] \mu_L) \quad (40)$$

which can be simplified as

$$\Pi = A \left(\underline{u} - c_0 \frac{(\underline{u} - \mu_L)(\mu_H - \underline{u})}{\lambda_{LH}} e^{\eta_0 R_0} \right). \quad (41)$$

Profit also depends on the inventory carrying costs. Let Ψ_i be defined as

$$\Psi_i = \int_{R_{i+1}}^{R_i} x (f_H(x) + f_L(x)) dx = c_i \frac{\Delta_i^L - \Delta_i^H}{\Delta_i^L} Q_i \quad (42)$$

where

$$Q_i = \begin{cases} ((\eta_i R_i - 1)e^{\eta_i R_i} - (\eta_i R_{i+1} - 1)e^{\eta_i R_{i+1}})/\eta_i^2 & \text{if } \eta_i \neq 0, \\ (R_i^2 - R_{i+1}^2)/2 & \text{if } \eta_i = 0. \end{cases}$$

The average inventory level **WIP** is

$$\mathbf{WIP} = \Psi_0 I_{\{Z(L)>0\}} + R_0 P^0 \quad (43)$$

where $I_{\{Z(L)>0\}}$ is an indicator function which is 1 if $Z(L) > 0$ and 0 otherwise.

Finally, the average profit per unit time in (2) is

$$\bar{J} = \Pi - g^+ \mathbf{WIP}. \quad (44)$$

The optimal values of $Z(L)$ and $Z(H)$ are determined by maximizing \bar{J} .

5.3 Other Performance Measures

We can also evaluate other measures related to the performance of the system. The average *sales rate* or *throughput rate* is

$$\mathbf{TH} = \lim_{T \rightarrow \infty} E \left[\frac{1}{T} \int_0^T u(\tau) d\tau \right] = \underline{u} - c_0 \frac{(\underline{u} - \mu_L)(\mu_H - \underline{u})}{\lambda_{LH}} e^{\eta_0 R_0}.$$

In addition to the throughput rate, various performance measures related to the service provided to the customers are also of interest. We measure the service by evaluating the ratio of the average sales to the average demand, by examining the fraction of customers who receive immediate service, and by the average number of waiting customers (the backlog level). The *service level*, the ratio of the average sales to the average demand, is

$$\mathbf{SL} = \mathbf{TH}/Ed.$$

The *fill rate* is the probability that a customer receives his product as soon as he arrives:

$$\mathbf{FR} = \lim_{t \rightarrow \infty} \mathbf{prob}[x(t) \geq 0] = P^0 + (P_0^H + P_0^L) I_{\{Z(L)>0\}}.$$

The *average backlog level* **BL** is

$$\mathbf{BL} = -\Psi_0 I_{\{Z(L)=0\}} - \sum_{j=1}^{J-1} \Psi_j - R_J P^J$$

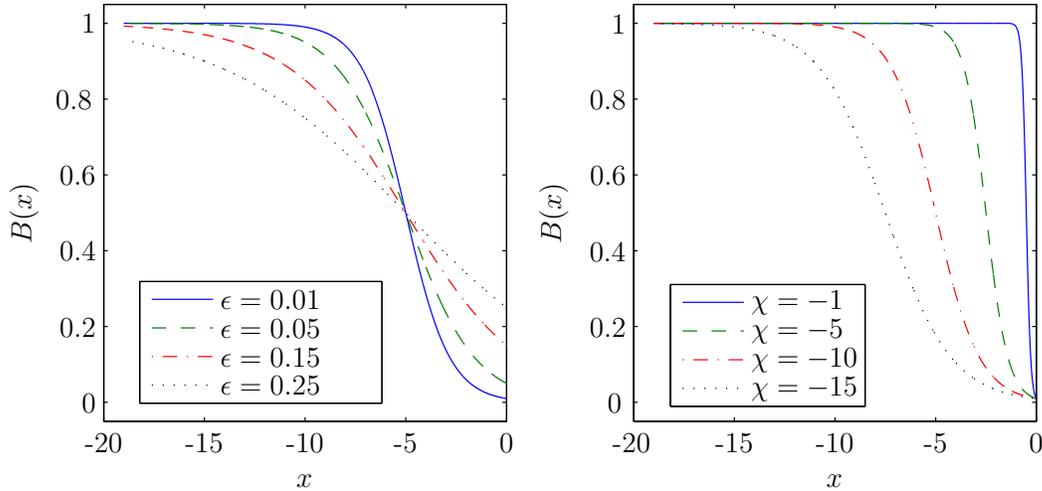


Figure 3: Various $B(x)$ functions for different values of ϵ when $\chi = -10$ and for different values of χ when $\epsilon = 0.01$

6 Behavior of the Model

Parametrization In the numerical examples of this section, $B(x)$ is a sigmoid function of the form

$$B(x) = \frac{1}{1 + e^{\gamma(x-\eta)}}. \quad (45)$$

where η is the median backlog tolerance, i.e., $B(\eta) = 0.5$, and γ determines the steepness of $B(x)$. Figure 1 shows a sigmoid function with $\gamma = 1/2$ and $\eta = -10$. While other functions satisfying (1) could be studied, our numerical results suggest that there is limited sensitivity to the exact shape of $B(x)$. Furthermore, it is shown in Veatch (2007) that the impact of $B(x)$ on system performance, for a given demand and production model, can be expressed through a single parameter. Thus, simple parametric forms of $B(x)$ are appropriate. Since our methodology is based on a piece-wise constant approximation of a given $B(x)$ function, the exact form of this function does not impose any limitations for our methodology.

In order to study the customer defection behavior, we focus on two parameters χ and ϵ where $\chi = 2\eta$ and $B(\chi) = 1 - \epsilon$. If ϵ is very small then χ can be considered as the maximum backlog after which almost all customers choose not to order. For example, setting $\chi = -10$ and $\epsilon = 0.01$ corresponds to a case where 99% of the potential customers choose not to order when the backlog level is -10 . With this parametrization, γ in Equation (45) is determined by χ and ϵ as $\gamma = \frac{2}{\chi} \ln\left(\frac{\epsilon}{1-\epsilon}\right)$. Consequently, for a given median backlog tolerance, ϵ also determines the steepness of $B(x)$. Figure 3 shows various $B(x)$ functions for different values of χ and ϵ .

In order to approximate a sigmoid defection function as a piece-wise constant function as in

Equation (9), we set

$$\begin{aligned} \beta_i &= -\delta i, & i = 0, 1, \dots, M+1 \\ B_i &= \frac{1}{2} \left(\frac{1}{1 + e^{\gamma(\beta_{i-1} - \eta)}} + \frac{1}{1 + e^{\gamma(\beta_i - \eta)}} \right), & i = 1, 2, \dots, M \\ B_{M+1} &= 1 \end{aligned}$$

where $\delta = \frac{1}{M\gamma} \ln \left(\frac{\epsilon'}{1-\epsilon'} \right) + \frac{\eta}{M}$ and ϵ' is a small positive number. With this setting, $B(x)$ reaches $1 - \epsilon'$ at the M th step and then reaches 1 at the $(M+1)$ st step and $\beta_{i+1} - \beta_i = \delta$ for $i = 0, 1, \dots, M$. In the numerical experiments, we used $M = 50$ and $\epsilon' = 10^{-4}$.

Effect of Customer Defection Behavior The effect of customer defection behavior on the performance of the system is examined by considering the effects of the backlog tolerance of the customers on the performance of the system (Figure 4). Figure 4 indicates that as customers become more sensitive to backlog, i.e., as χ increases, the profit and the service level decrease, and the expected inventory level increases. Due to the loss of more and more customers, the expected backlog level also decreases and the service level is low. The upper and lower boundaries (Z_L and \underline{X}) increase as the customers become more sensitive to the waiting time. Note that $\chi = 0$ is the lost sales case and $\chi \rightarrow -\infty$ is the complete backordering case. Most of the benefit of customer patience is already seen at $\chi = -10$.

Figure 5 shows that the effect of the steepness of the customer defection function when the median delay tolerance is fixed. Figure 5 shows that the effect of ϵ on the system performance is limited compared to the effect of χ especially when $\epsilon \ll 1$. This suggests that the median backlog tolerance, $\eta = \frac{\chi}{2}$ of customers can be inferred from the managers to form the $B(x)$ function by setting ϵ to a predetermined low value such as 0.01.

Effect of Customer Variability In order to capture the effect of the variability of demand on the performance of the system, consider $N(t)$ which has been defined as the cumulative demand during $[0, t]$ (Section 2.3). The variance of $N(t)$ per unit time in the long run can be written as

$$\lim_{t \rightarrow \infty} \frac{\text{Var}[N(t)]}{t} = \frac{2(\mu_H - \mu_L)^2(1 - e)}{\lambda_{LH}}$$

(Tan 1997). Therefore, for fixed μ_H , μ_L , and average demand rate, the variability of demand depends on λ_{LH} .

Figure 6 shows the effect of λ_{LH} on the system performance. In this example, since the machine has sufficient capacity to satisfy the demand in the long run, as the demand variability decreases, i.e. when λ_{LH} increases, the hedging point $Z(L)$ also decreases. Similarly, the average backlog and the average inventory decrease and the profit and service level increase. Note that if there is no variability in the demand and the machine has sufficient capacity, then it would be possible to meet the demand with no inventory or backlog and the profit would reach $A\underline{u}$ (2.7 in the case shown in Figure 6).

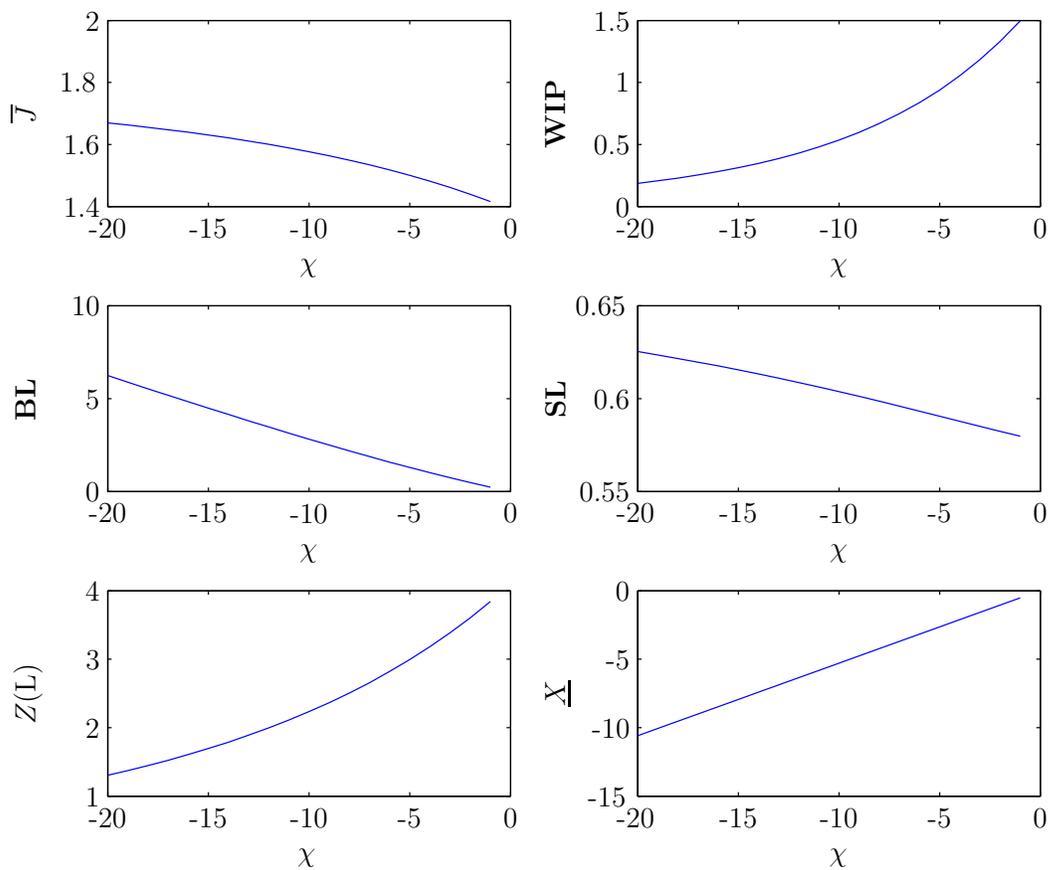


Figure 4: Effect of Median Delay Tolerance ($\mu_H = 1.5, \mu_L = 0.3, \lambda_{HL} = 0.05, \lambda_{LH} = 0.05, \underline{u} = 0.6, A = 3, g^+ = 0.1, \epsilon = 0.01$)

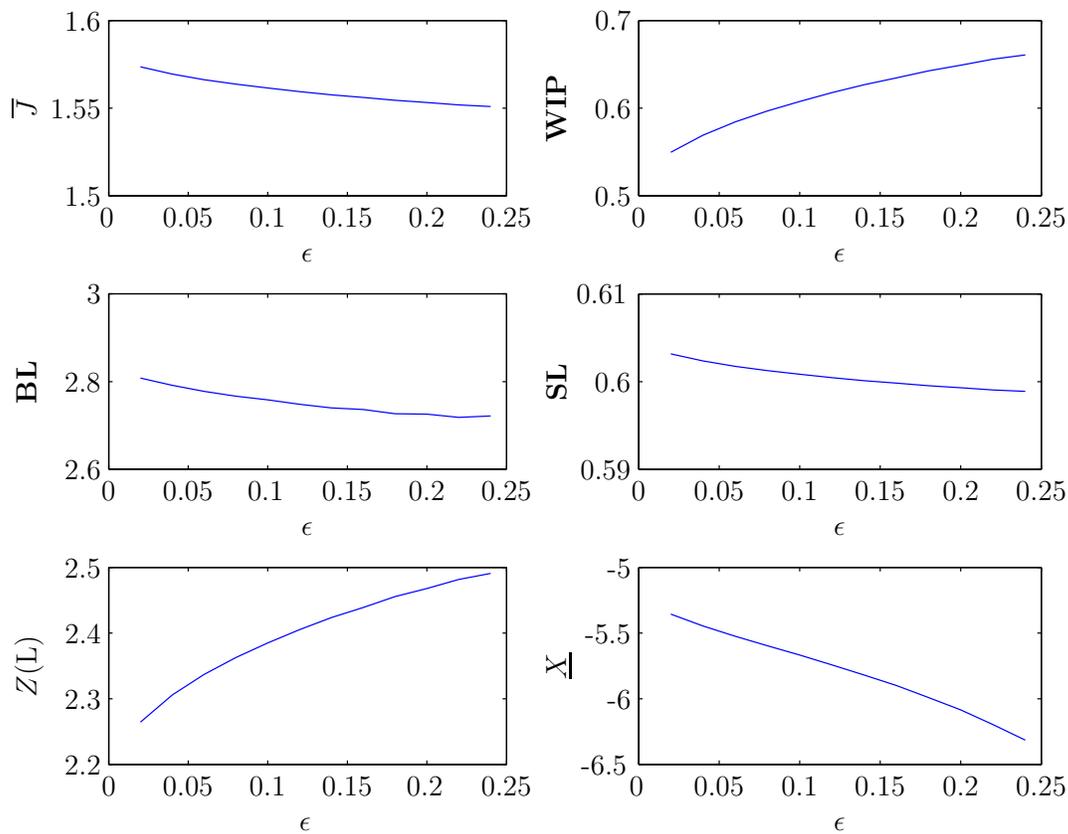


Figure 5: Effect of Delay Tolerance Steepness ($\mu_H = 1.5, \mu_L = 0.3, \lambda_{HL} = 0.05, \lambda_{LH} = 0.05, \underline{u} = 0.6, A = 3, g^+ = 0.1, \chi = -10$)

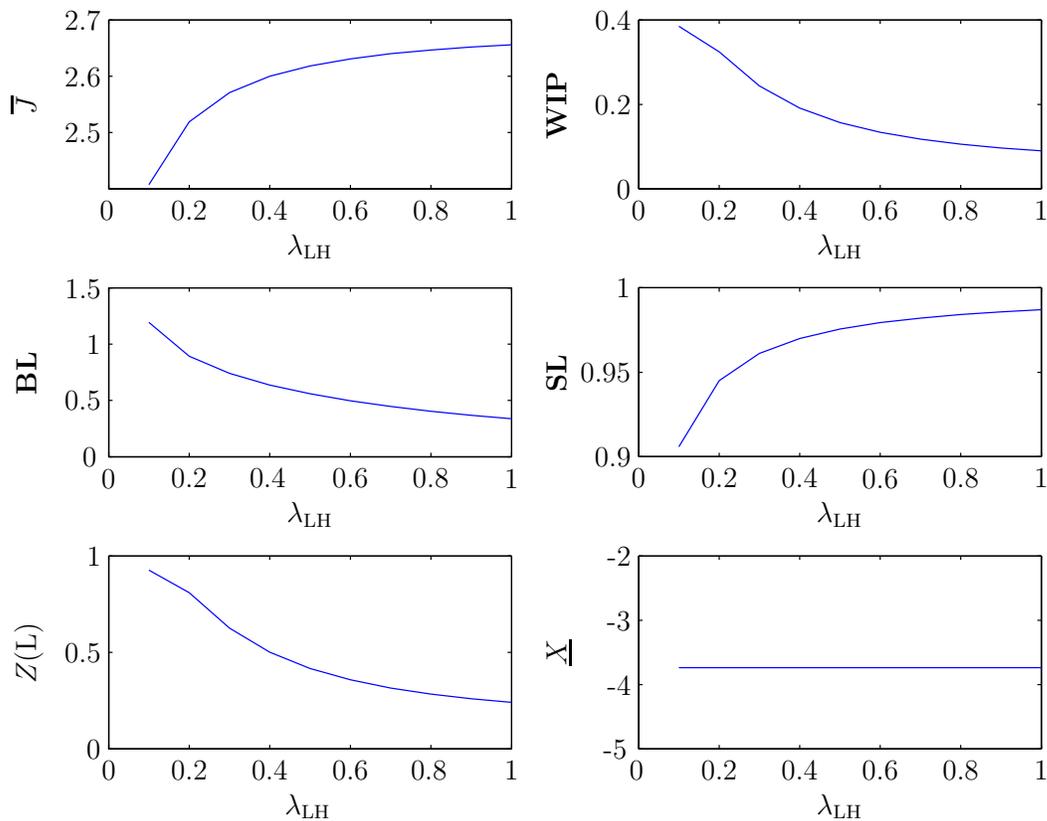


Figure 6: Effect of the demand variability ($\mu_H = 1.5, \mu_L = 0.3, Ed = 0.9, \underline{u} = 1.1, \lambda_{HL} = \frac{\mu_H - Ed}{Ed - \mu_L} \lambda_{LH}, A = 3, g^+ = 0.1, \chi = -10, \epsilon = 0.01,)$

7 Conclusions

We have extended a widely-studied dynamic programming model of real-time scheduling control of manufacturing systems in two important ways: we model the effect of backlog on profits through an explicit representation of customer behavior; and we model random demand.

The new model of customer behavior involves a *defection function* which indicates what fraction of the potential customers choose not to complete their orders when the backlog reaches a given level. Because of this phenomenon, the model has a novel feature: the demand need not be less than capacity for there to exist a steady-state probability distribution of the inventory/backlog and the demand state.

We use a relationship with the lost sales problem to determine a solution structure, and we find that the solution involves a hedging point (to limit how far production should be allowed to go ahead of demand). To determine the hedging point, we find the steady-state probability distribution. We evaluate the objective function and choose values of the hedging point to maximize it.

Finally, we have performed a set of numerical experiments to demonstrate the behavior of the new model and the solution.

Acknowledgments

We are grateful for support from the TUBITAK and the TUBA, the National Science Foundation, Grant DMI-9713500, the Lean Aircraft Initiative, the Xerox Foundation, and the Singapore-MIT Alliance. We are also grateful for the suggestions provided by the reviewers and the department editor.

Appendix—Proof of Theorem 1

First we show that the optimal policy is non-idling when there are backorders.

Lemma 1 $u^*(x, D) = \underline{u}$ for $x < 0$.

Proof. Suppose $u^*(x, D) < \underline{u}$ on a subset of $[x^*, 0]$ of nonzero measure. Consider the coupled processes $x(t)$, $x^1(t)$ with policies u^* and u^1 and identical initial conditions with $x(0) < 0$. Let $u^1(x, D) = \underline{u}$ for $x < 0$, $u^1(0, L) = \mu_L$ for $t < \tau \equiv \min\{t : x(t) = 0\}$ so that $x^1(t) \leq 0$ until the processes merge at $x(\tau) = x^1(\tau) = 0$, and $u^1(t) = u^*(t)$ for $t > \tau$. Then

$$E \int_0^\tau d(t)[1 - B(x(t))]dt < E \int_0^\tau d(t)[1 - B(x^1(t))]dt,$$

i.e., $x(t)$ experiences less demand, in expectation. The processes merge with probability one and the only cost difference is due to the difference in demand before merging:

$$\lim_{T \rightarrow \infty} J^1(x, D; T) - J(x, D; T) = -E \int_0^\tau Ad(t)[B(x^1(t)) - B(x(t))]dt > 0,$$

contradicting the optimality of u^* . ■

Proof of Theorem 1

Let $Z(L) = \min\{x : u^*(x, L) < \underline{u}\}$. Because (7) is linear in u , we can choose $u^*(\cdot, L)$ to only have the values 0 and \underline{u} , except at points where both maximize (7), where we can choose the value μ_L . Consequently, $Z(L)$ is an upper bound: if $x(0) \leq Z(L)$, then $x(t) \leq Z(L)$.

If $Z(L) = 0$ we are done, so we can assume $Z(L) > 0$. Temporarily let x denote a random variable with the stationary distribution of $x(t)$ and similarly for $d(t)$ and $B(x(t))$. Equation (44) can be written

$$\bar{J} = A(Ed - \mathbf{prob}[x \leq 0]\mu_B) - g^+ \mathbf{prob}[x > 0]E(x|x > 0) \quad (46)$$

where $\mu_B = E(dB|x \leq 0)$ is the average rate at which customers balk when there is a backlog, and the stationary distribution is used. Because entry into $x \leq 0$ occurs at a single state $(0, H)$, the conditional distribution of (x, D) given $x \leq 0$ is the same for all policies satisfying Lemma 1. In particular, μ_B is the same.

We will establish a relationship between the backlog-dependent problem $x(t)$ and the lost sales problem $x^{LS}(t)$ of Hu (1995). The lost sales problem has the constraint $x^{LS}(t) \geq 0$ (no backlogging) and machine failures in stead of random demand. Its dynamics can be written

$$\frac{dx^{LS}(t)}{dt} = u(t) + \underline{u} - \mu_H, \quad \begin{array}{l} 0 \leq u(t) \leq \mu_H - \mu_L, D = L \\ u(t) = 0, D = H \end{array} \quad (47)$$

where we have identified the operational state as L, the failed state as H, the demand rate $\mu_H - \underline{u}$, and the production rate $\mu_H - \mu_L$. On $x > 0$, (3) - (5) can be written

$$\frac{dx(t)}{dt} = u(x(t)) + \underline{u} - \mu_H, \quad \begin{array}{l} (\mu_H - \mu_L) - \underline{u} \leq u(t) \leq \mu_H - \mu_L, D = L \\ -\underline{u} \leq u(t) \leq 0, D = H \end{array} \quad (48)$$

Thus, on $x > 0$ our model is a relaxation of $x^{LS}(t)$ in which the lower control limits are negative. Temporarily let x^{LS} denote a random variable with the stationary distribution of $x^{LS}(t)$. Average cost for the lost sales problem is

$$\bar{J}^{LS} = s \mathbf{prob}[x^{LS} = 0] + g^+ E(x^{LS}|x^{LS} > 0) \mathbf{prob}[x^{LS} > 0], \quad (49)$$

where s is the stockout cost. Let $Z(L)^{LS}$ be the optimal hedging point for $x^{LS}(t)$, and also assume $Z(L)^{LS} > 0$, so that $\mathbf{prob}[x^{LS} = 0, D = L] = 0$ and (49) is equivalent to applying a cost to lost sales.

Suppose $x(t)$ and $x^{LS}(t)$ both use the optimal policy for $x^{LS}(t)$ on $[0, \infty]$ and that $x(t)$ uses the optimal policy of Lemma 1 on $(-\infty, 0)$. To compare their differential costs, set $V^u(0, H) = V_{LS}^u(0, H) = 0$. For some s , $V^u(0, L) = -V_{LS}^u(0, L)$, because V_{LS}^u is a continuous function of s (see also (16) of (Hu 1995)). Then, because they have the same dynamics on $x > 0$, entry into $x > 0$ occurs from a single state $(0, L)$, and departure from $x > 0$ occurs at the recurrent state $(0, H)$, their differential costs on $x > 0$ are the same:

$$V^u(x, D) = -V_{LS}^u(x, D), \quad x > 0. \quad (50)$$

The negative sign is needed because $x^{\text{LS}}(t)$ is a minimization problem. Hu (1995) finds V_{LS} and shows that $V_{\text{LS}}(x, L)$ is decreasing for $x < Z(L)^{\text{LS}}$ and increasing for $x > Z(L)^{\text{LS}}$. It can also be verified that $V_{\text{LS}}(x, H)$ is decreasing for $x < Z(H)^{\text{LS}}$, where $Z(L)^{\text{LS}} < Z(H)^{\text{LS}}$. Thus, $V^u(x, D)$ is increasing and satisfies the Bellman equations (7) and (8) for $0 < x < Z(L)^{\text{LS}}$. By Lemma 1, V^u can be extended to also satisfy these equations for $x \leq 0$. Since V^u is fully specified, it must be the unique optimal differential cost, $V(x, D) = V^u(x, D)$, for $x < Z(L)^{\text{LS}}$.

Finally, V is continuously differentiable (see, e.g., Sethi, Suo, Taksar, and Zhang 1997) so $\frac{dV}{dt}(Z(L)^{\text{LS}}, L) = \frac{dV_{\text{LS}}}{dt}(Z(L)^{\text{LS}}, L) = 0$, and we can choose $u^*(Z(L)^{\text{LS}}, L) = \mu_L$. That makes $Z(L) = Z(L)^{\text{LS}} > 0$ a hedging point for both problems and states $x > Z(L)$ transient. Hence, the Bellman equation is satisfied in all recurrent states by the policy of Theorem 1, and it is optimal. ■

Remark. Another approach to Theorem 1 is to extend proofs of convexity of the differential cost, such as Sethi, Suo, Taksar, and Zhang (1997), to show that V is concave. Concavity implies that the hedging point form extends to all states. However, the connection we have made with the lost sales problem makes explicit expressions for V available. The connection with the lost sales problem (or any problem with the same dynamics and cost rate on $x > 0$) can also be seen by comparing (46) with (49) and recalling that μ_B is the same for all policies satisfying Lemma 1. Of all policies with a given **prob** $[x > 0]$, the optimal policy minimizes $E(x|x > 0)$. In (49), again the optimal policy minimizes $E(x^{\text{LS}}|x^{\text{LS}} > 0)$ for a given **prob** $[x^{\text{LS}} > 0]$. When they use the same policy on the recurrent states $x \leq Z(L)$, they have $E(x|x > 0) = E(x^{\text{LS}}|x^{\text{LS}} > 0)$. Thus, policies have the same average cost ordering for both problems. This similarity demonstrates, as noted above, that our problem differs from the lost sales problem only in that the control limits are relaxed.

References

- Anderson, E., G. J. Fitzsimons, and D. Simester (2003). Mitigating the cost of stockouts. Technical report, University of Chicago.
- Bielecki, T. and P. R. Kumar (1988). Optimality of zero-inventory policies for unreliable manufacturing systems. *Operations Research* 36(4), 532–541.
- Brown, L., N. Gans, A. Mandelbaum, A. Sakov, H. Shen, S. Zeltyn, and L. Zhao (2003). Statistical analysis of a telephone call center: A queueing-science perspective. Center for Financial Institutions Working Papers 03-12, Wharton School Center for Financial Institutions, University of Pennsylvania. available at <http://ideas.repec.org/p/wop/pennin/03-12.html>.
- Dellaert, B. G. C. and B. E. Kahn (1999). How tolerable is delay?: Consumers' evaluations of internet web sites after waiting. *Journal of Interactive Marketing* 13(1), 41–54.
- Fitzsimons, G. J. (2000). Consumer response to stockouts. *Journal of Consumer Research* 27, 249–266.
- Fleming, W. H., S. P. Sethi, and H. M. Soner (1987). An optimal stochastic production planning with randomly fluctuating demand. *SIAM Journal of Control Optimization* 25, 1495–1502.
- Ghosh, M. K., A. Araposthathis, and S. I. Markus (1993). Optimal control of switching diffusions with applications to flexible manufacturing systems. *SIAM Journal of Control Optimiza-*

- tion 31, 1183–1204.
- Gross, D. and C. M. Harris (1985). *Fundamentals of Queueing Theory (2nd edition)*. John Wiley & Sons, Inc., New York, NY.
- Hall, R. W. (1991). *Queueing Methods for Services and Manufacturing*. Englewood Cliffs, NJ: Prentice Hall.
- Hu, J. (1995). Production rate control for failure prone production with no backlog permitted. *IEEE Transactions on Automatic Control* 40(2), 291–295.
- Kimemia, J. G. and S. B. Gershwin (1983). An algorithm for the computer control of production in a flexible manufacturing systems. *IIE Transactions* 15(4), 353–362. Reprinted in *Modeling and Control of Automated Manufacturing Systems*, ed. Alan A. Desrochers, IEEE Computer Society Press Tutorial, 1990.
- Mandelbaum, A. and S. Zeltyn (2004). The impact of customers patience on delay and abandonment: some empirically-driven experiments with the m/m/n+g queue. *OR Spectrum* 26, 377411.
- Martin, S. and P. C. Smith (1999). Rationing by waiting lists: An empirical investigation. *Journal of Public Economics* 71(1), 141–164.
- Martinelli, F. and P. Valigi (2004). Hedging point policies remain optimal under limited backlog and inventory space. *Automatic Control, IEEE Transactions on* 49(10), 1863–1871.
- Olsder, G. J. and R. Suri (1980, December). Time-optimal control of flexible manufacturing systems with failure prone machines. In *Proceedings of the 19th IEEE Conference on Decision and Control*, Albuquerque, New Mexico.
- Perkins, J. and R. Srikant (2001). Failure-prone production systems with uncertain demand. *Automatic Control, IEEE Transactions on* 46(3), 441–449.
- Raiffa, H. (1982). *Art and Science of Negotiation*. Harvard University Press.
- Rishel, R. (1975, February). Dynamic programming and minimum principles for systems with jump markov disturbances. *SIAM Journal on Control* 13(2), 338–371.
- Sethi, S., W. Suo, M. Taksar, and Q. Zhang (1997). Optimal production planning in a stochastic manufacturing system with long-run average cost. *Journal of Optimization Theory and Applications* 92, 161–188.
- Tan, B. (1997). Variance of the throughput of an N -station production line with no intermediate buffers and time dependent failures. *European Journal of Operational Research* 101(3), 560–576.
- Tan, B. (2002). Production control of a pull system with production and demand uncertainty. *IEEE Transactions on Automatic Control* 47(5), 779–783.
- Tan, B. and S. B. Gershwin (2001). On production and subcontracting strategies for manufacturers with limited capacity and volatile demand. Working Paper Series ORC 354-01, Massachusetts Institute of Technology, Operations Research Center.

- Taylor, S. (1994). Waiting for service: The relationship between delays and evaluations of service. *Journal of Marketing* 58(2), 56–69.
- Veatch, M. (2007). The impact of customer impatience on production control. *IIE Transactions*. To appear.
- Weinberg, B. D. (2000). Don't keep your internet customers waiting too long at the (virtual) front door. *Journal of Interactive Marketing* 14(1), 30–39.
- Yeralan, S. and B. Tan (1997). A station model for continuous materials flow production. *International Journal of Production Research* 35(9), 2525–2541.