

**OPTIMIZING DEMAND MANAGEMENT IN STOCHASTIC  
SYSTEMS TO IMPROVE FLEXIBILITY AND PERFORMANCE**

A Thesis  
Presented to  
The Academic Faculty

by

Serhan Duran

In Partial Fulfillment  
of the Requirements for the Degree  
Doctor of Philosophy in the  
H. Milton Stewart School of Industrial and Systems Engineering

Georgia Institute of Technology  
August 2007

**OPTIMIZING DEMAND MANAGEMENT IN STOCHASTIC  
SYSTEMS TO IMPROVE FLEXIBILITY AND PERFORMANCE**

Approved by:

Professor Julie L. Swann, Advisor  
H. Milton Stewart School of Industrial  
and Systems Engineering  
*Georgia Institute of Technology*

Professor Hayriye Ayhan  
H. Milton Stewart School of Industrial  
and Systems Engineering  
*Georgia Institute of Technology*

Professor Mark Ferguson  
School of Management  
*Georgia Institute of Technology*

Professor Paul M. Griffin  
H. Milton Stewart School of Industrial  
and Systems Engineering  
*Georgia Institute of Technology*

Professor Pınar Keskinocak  
H. Milton Stewart School of Industrial  
and Systems Engineering  
*Georgia Institute of Technology*

Date Approved: 11 June 2007

*To Mom and Dad,*

*Nursel and Kayhan Duran,*

*for their unconditional love and support.*

## ACKNOWLEDGEMENTS

My sincere thanks belong to my advisor Dr. Julie L. Swann. Her contributions on molding of my research perspective will have an everlasting impact. Her knowledge, experience and insights have been very influential and helpful in my studies. I would also like to express my special thanks to Dr. Pınar Keskinocak, without her guidance the accomplishment of this thesis would not be possible. I would like to thank Dr. Hayriye Ayhan, Dr. Mark Ferguson, and Dr. Paul M. Griffin. I appreciate their time and effort in serving on my Ph.D. committee.

## TABLE OF CONTENTS

DEDICATION . . . . .	iii
ACKNOWLEDGEMENTS . . . . .	iv
LIST OF TABLES . . . . .	viii
LIST OF FIGURES . . . . .	ix
SUMMARY . . . . .	x
I INTRODUCTION . . . . .	1
II OPTIMAL PRODUCTION AND INVENTORY POLICIES OF PRIORITY AND PRICE-DIFFERENTIATED CUSTOMERS . . . . .	5
2.1 Introduction . . . . .	5
2.1.1 Motivation and Background . . . . .	5
2.1.2 Literature Review . . . . .	7
2.2 Models and Results . . . . .	9
2.2.1 Notation and Assumptions . . . . .	10
2.2.2 Priority Differentiation Strategy (PDS) . . . . .	12
2.2.3 Special Cases and Extensions . . . . .	19
2.3 Computational Analysis . . . . .	21
2.3.1 Experiment Details . . . . .	21
2.3.2 Results . . . . .	23
2.4 Conclusions . . . . .	25
III POLICIES UTILIZING TACTICAL INVENTORY FOR SERVICE DIFFEREN- TIATED CUSTOMERS . . . . .	31
3.1 Introduction . . . . .	31
3.2 Models and Results . . . . .	33
3.2.1 Assumptions and Notation . . . . .	33
3.2.2 Time Differentiation Strategy . . . . .	35
3.2.3 Common Service Strategy . . . . .	37
3.3 Results . . . . .	39
3.4 Computational Analysis . . . . .	42

3.4.1	Experiment Details . . . . .	43
3.4.2	Results . . . . .	43
3.5	Conclusions . . . . .	44
IV	LEADTIME QUOTATION AND ORDER ACCEPTANCE WHEN DEMAND DEPENDS ON SERVICE PERFORMANCE . . . . .	47
4.1	Introduction . . . . .	47
4.2	Literature Review . . . . .	48
4.3	Models and Results . . . . .	52
4.3.1	Infinite Capacity Case . . . . .	52
4.3.2	Finite Capacity Case . . . . .	57
4.4	Numerical Experiments . . . . .	63
4.4.1	Experiment Details . . . . .	63
4.4.2	Results . . . . .	64
4.5	Conclusions . . . . .	68
V	DYNAMIC SWITCHING TIMES FOR SEASON AND SINGLE TICKETS IN SPORTS AND ENTERTAINMENT . . . . .	70
5.1	Introduction . . . . .	70
5.2	Literature Review . . . . .	72
5.3	Assumptions and Notation . . . . .	74
5.4	Model and Results . . . . .	75
5.5	Extensions . . . . .	80
5.5.1	$\ell$ -Performances ( $\ell > 2$ ) During the Selling Period . . . . .	80
5.5.2	Time-Dependent Demand Rates . . . . .	81
5.6	Computations . . . . .	82
5.7	Conclusions . . . . .	87
VI	CONCLUSION . . . . .	89
APPENDIX A	PROOFS FOR CHAPTER 2 . . . . .	92
APPENDIX B	PROOFS FOR CHAPTER 3 . . . . .	102
APPENDIX C	PROOFS FOR CHAPTER 4 . . . . .	114
APPENDIX D	PROOFS FOR CHAPTER 5 . . . . .	121

REFERENCES . . . . .	126
VITA . . . . .	132

## LIST OF TABLES

1	Notation . . . . .	12
2	Additional notation . . . . .	13
3	Specific experimental data . . . . .	22
4	Additional notation . . . . .	35
5	Specific experimental data . . . . .	44
6	Optimal accept-up-to level when the service level is 0 . . . . .	60
7	Optimal accept-up-to level when there are six customers in the system . . .	61
8	Selected Optimal Switching Times when $p_B = 220$ . . . . .	84



## LIST OF FIGURES

1	Optimal Decisions under the Optimal Policies for PDS . . . . .	27
2	The relative performance of $(S, R, B)$ type policies over traditional . . . . .	28
3	Tactical Inventory Levels with Price Increasing over the Horizon . . . . .	29
4	The relative performance of PDS over traditional with different prices . . . . .	30
5	Optimal Decisions for TDS under the Optimal Policies . . . . .	42
6	The relative performance of $(S, R, B)$ type policies over traditional inventory decisions . . . . .	46
7	Performance based on revenue per item when $\alpha = 0.2, \mu = 0.5, c = 2$ . . . . .	55
8	Performance based on penalty cost when $\alpha = 0.5, \mu = 0.5, R = 2$ . . . . .	56
9	Percentage Improvement of Optimal over Naive with the Horizon Length . . . . .	64
10	Percentage Improvement of Optimal over Naive with the Number of Service Levels . . . . .	65
11	Percentage Improvement of Optimal over Naive with Service rate . . . . .	66
12	Percentage Improvements of Optimal and Heuristic over Naive with Increasing Cost . . . . .	67
13	Optimal Switching Times at Different Bundle Prices . . . . .	85
14	Optimal Switching Times for Various Parameters . . . . .	86
15	Comparison of Dynamic Decision of Switch Time vs Static . . . . .	87

## SUMMARY

In this thesis we analyze optimal demand management policies for stochastic systems. In the first system considered, a manufacturer decides how to manage demand from customers that differ in their priority level and willingness to pay. He has limited production capacity and predetermined prices throughout the horizon. We find an optimal production and inventory strategy that rations current and future limited capacity between customer classes through reserving inventory for the future and accepting orders now for future delivery. Next, we extend these results to the case when the customers have different tolerance to delayed fulfillment, namely, first-class customers never accept backlogging whereas second-class customers agree to wait one period for a discount. We find an optimal policy similar to the production and inventory strategy that is used for the first system based on threshold values. The third system considers a firm whose recent performance in meeting quoted leadtimes affects future demand arrivals. We assume that the probability of a customer placing an order depends on the quoted leadtime, and both customer arrivals and processing times are stochastic. When capacity of the firm is infinite, we find the optimal leadtime to quote, and when capacity is finite and leadtime is industry-dictated, we determine that the optimal demand acceptance policy does not necessarily have a nice structure. We comment on the structure of the optimal policy for a special case and develop several heuristics for the general case. The final system considered in this thesis is the Sports and Entertainment industry, where demand is managed for a season of several performances by selling season tickets initially and single events later in the selling horizon. We specifically study the optimal time to switch between these market segments dynamically as a function of the state of the system and show that the optimal switching time is a set of time thresholds that depend on the remaining inventory and time left in the horizon.

# CHAPTER I

## INTRODUCTION

Matching supply and demand effectively is one of the key factors for a profitable business. The ill-management of supply can cause excess production, inventory and labor costs or loss of potential revenue, whereas incompetent demand management policies may result in the loss of customers and their loyalty. Research on demand management has been pursued extensively by scholars (economists and operations researchers) due to its importance and potential for improvement in system. However, rapid advances in technology and science bring shifts in the business environment in which companies are operating, and these shifts require new models and techniques to be utilized, therefore continued efforts are required to improve study in this area. Pricing, rationing and product/service differentiation policies are several tools to manage demand effectively. In this thesis, we focus on demand management under stochastic operating environments using tools other than pricing.

Demand management is usually considered as a separate activity from production operations, and in practice most companies have a marketing department deciding pricing, promotion and advertisement policies to manipulate demand without consulting the production department or considering capacity limitations. Intuitively, performance improvements using demand management that considers production operations or capacity limitations are expected to be higher, and empirical evidence confirms this expectation (see Hausman et al. [42]). By this motivation, this dissertation aims at optimizing demand management in manufacturing and service systems under stochastic operating environments considering production and capacity limitations. The goal is to use demand management to increase flexibility to the firm, which can increase profits and improve customer service.

The first part of this dissertation (Chapter 2) addresses demand management of customers who differ in their priority level and willingness to pay, by a firm with limited

production capacity. Today, many firms are exploring production and supply chain strategies where customers may be segmented into different classes based on service level or priority, which can result in a more efficient production system as well as a better match between supply and demand. Specifically, when customers are segmented into classes by service levels based on delivery time, customers with an immediate need (e.g., businesses) receive expedited product, while flexible customers receive incentives for their patience. The firm benefits from the flexibility in production gained by backlogging or from longer leadtime requirements, enabling it to meet more demand or use less overtime to satisfy the same number of customers. An example of a company using differentiation is Amazon.com, where consumers can choose expedited shipping or free shipping. In the latter Amazon.com receives increased flexibility, since the stated leadtime exceeds the actual processing and transportation time.

Specifically, in Chapter 2 we focus on production and inventory decisions of a firm using stochastic inventory control, operating in an environment where customer classes are differentiated by their priority level. We introduce “tactical inventory decisions” to improve the profit, service, and flexibility of the system. The tactical decisions include the use of inventory or capacity allocations in one time period to serve customer demand in another time period. Specifically, we allow the firm to reserve inventory to satisfy future demand and to plan backlogging to serve current demand. We analyze the structure of optimal production and inventory strategies that result when customer classes differ in priority. We find that a set of threshold policies for the production, reserving and backlogging decisions is optimal even with multiple classes, and the policy is nested by class. We perform computational analysis to see that the profit attained when strategic inventory is used can be a significant improvement over a traditional inventory policy.

In Chapter 2, customer classes are assumed to be differentiated by their priority level, where higher-class customers receive complete priority over the lower-class customers in the use of current resources and future backlogging. A key feature of this work is that customers in both of the classes behave homogeneously in terms of the delivery time, i.e., all customers are willing to wait for fulfillment. But in practice, the segmented customer

classes may behave differently with respect to their acceptance of delay fulfillment, i.e., one class may never accept delay fulfillment. Therefore in Chapter 3, we develop models that incorporate tactical inventory decisions for customer classes with non-homogeneous behavior in terms of delay fulfillment. Specifically, we assume that the first-class customers claim the item immediately and never accept to be backlogged, whereas second-class customers accept delay fulfillment for a discount. Although the proof techniques are similar to the ones in Chapter 2, the results do not follow directly since the difference in customers' tolerance to delay fulfillment changes the structure of the models. We show that the optimal production and inventory strategies for patient and impatient customers are threshold policies for the production, reserving and backlogging decisions, as the ones that are proven to be optimal in Chapter 2 for priority-differentiated customers.

In the third part of this dissertation (Chapter 4), we consider the optimal demand management of a firm when customers choose the firm according to a firm's past performance. Mostly thanks to the Internet, in the current business environment customers can easily share their experiences with each other, informing customers' decisions about whether or not to do business with a firm or buy that firm's products. The delivery time and price of an item are not the only factors that affect the customer's decision, but customers may also consider the past performance of a producer, specifically whether he is meeting the promised delivery times or not. For example, 78% of companies that operate in a just-in-time environment in the U.S. ranked delivery reliability as high priority, whereas only 25% ranked price as high priority (Billesbach et al. [9]).

Specifically, we consider the optimal demand management of a firm via leadtime quotation and order acceptance when the firm's recent performance of meeting quoted leadtimes impacts future orders from customers. For this research, allowing leadtime performance to impact future customer arrivals is an idea that we introduce into the model, since this may be true in practice. We consider the problem for both infinite and finite capacity cases. For the infinite capacity case, we find the optimal closed-form expression for leadtime quotation. We show that the optimal leadtime to quote that accounts for past performance is more conservative (i.e., longer) than the optimal leadtime that ignores it. We also find

that the optimal leadtime is always positive, unlike in the case that ignores service, which means that a firm considering the past performance effect would never quote an *unethical* leadtime of zero. When capacity is finite and leadtime is industry-dictated, we determine that the optimal demand acceptance policy does not necessarily have a nice structure, but in some special cases it is convex in the service level of the firm. For the finite capacity case, we also develop several heuristics for the order acceptance model with general stochastic production.

In the final part of the dissertation, Chapter 5, we are analyzing demand management for the sports and entertainment industry via the selling of season tickets vs. single tickets. Common industry practice in the sports events is *pure bundling*, selling only season tickets first and switching to single ticket sales later in the selling horizon. We will address the issue of dynamically deciding when to switch from season tickets to singles by considering the optimal stopping time, which will enable us to take the actual sales realization into consideration.

Initially, we consider a two-performance selling season and the processes for the bundled and single-tickets to be Poisson processes with constant rates. These assumptions are later relaxed in the chapter. We show that the optimal time to switch is determined by a set of threshold pairs, which are defined by the remaining inventory and the time left in the horizon. After each sale, the current time is compared to the time threshold for the corresponding remaining inventory to determine if the switch should be made immediately or not. We also perform numerical experiments to illustrate the value of dynamically deciding the switching time instead of deciding it without observing any sales realization, and we report significant percentage improvements in revenue.

For each of the four topics in the thesis, we present a review of the literature in the corresponding chapter, describe how our work contributes to the literature, and present the main models and results. Major proofs are provided in the Appendix. We conclude the thesis by identifying several areas of future research.

## CHAPTER II

### OPTIMAL PRODUCTION AND INVENTORY POLICIES OF PRIORITY AND PRICE-DIFFERENTIATED CUSTOMERS

#### *2.1 Introduction*

##### **2.1.1 Motivation and Background**

Flexibility is essential for businesses in order to deal with variability, uncertainty, and changes in the business environment. Manufacturing flexibility can be achieved in many ways including labor force, machinery, product mix, product design, or new products. Increasingly, companies are also turning to customer segmentation and tactical inventory decisions as a source of flexibility.

Differentiated service levels based on delivery time allow customers with an immediate need (e.g., businesses) to receive expedited product, while flexible customers receive incentives for their patience. An example of a company using differentiation is Amazon.com, where consumers can choose expedited shipping or free shipping. In the latter Amazon.com receives increased flexibility, since the stated leadtime exceeds the actual processing and transportation time. Customer segmentation by time, whether in manufacturing or the airline industry, provides a mechanism for balancing the supply and demand requirements of the system (e.g., shifting leisure travel from Friday to Saturday), which allows more efficient use of existing resources. A key example of a manufacturing company that employs flexibility in managing customer demand is Dell Inc. Customers are segmented according to type (e.g., business versus personal), and prices of products change regularly [1].

The primary goal of this research is to provide tools for managing production and inventory tactically when customers differ in their willingness to pay and their willingness to wait. The key questions we address are how much to produce and how to allocate scarce resources (either current inventory or future limited production capacity) dynamically among different customer classes. We incorporate a firm's tactical inventory decisions, which

we define to mean inventory or capacity allocations in one time period to serve customer demand in another time period. Specifically, we allow the firm to reserve inventory to satisfy future demand (sometimes called “discretionary sales”), and to plan backlogging, where the firm can accept orders in a period to be delivered in the future.

For example, many manufacturing companies face the following problem: some customers are willing to pay high prices to receive faster fulfillment; other customers are willing to accept a lower priority for fulfillment, but they demand low prices. The manufacturer has limited production capacity, and in order to maximize the profit, he needs to allocate the capacity effectively. With an advanced strategy, the manufacturer can separate the customers into multiple classes according to priority levels and then manage the production and the inventory appropriately; we refer to this as a *differentiated* strategy.

In this chapter we study the *Priority Differentiation Strategy* (PDS), where we assume the first class pays a premium to have higher priority in the current period over production and inventory resources compared to the second class. We assume that the manufacturer can or is willing to prioritize demand classes. That is, the manufacturer makes a decision on higher priority demand before he accepts or rejects the lower priority demand requests. This situation might occur in practice when requests are submitted electronically and are handled in batches, or it could result from any working environment where a manufacturer may temporarily ignore requests from second-class customers. Studying the general model also allows us to analyze several situations that are special cases or extensions of it. For example, in some circumstances the manufacturer is not able or not allowed to differentiate the customers and will deal with them as a single class.

We assume demand in each period is a general function of price, is continuous and differentiable, and is lost if rejected; we do not make restrictive assumptions regarding the stochastic demand arrivals and the production process. We focus on a periodic review environment where prices are predetermined but not known by customers until the current period. We assume backordered demand is fulfilled in the next period.



### 2.1.2 Literature Review

One stream of literature related to our work is inventory theory, especially when there are multiple classes of customers. Two seminal papers in this area are Veinott [83] and Topkis [80]. In [83], Veinott shows some conditions under which a base-stock policy is optimal for the production decision when cost minimization is the goal. When parameters are time varying and the classes have different priorities, the demand from a higher class should be satisfied before demand from a lower class, and further restrictions are necessary on the costs. A related topic is considered in [80], where the work is extended to decide a set of critical levels that determine when to satisfy a particular class of demand. Topkis outlines some assumptions under which the optimal policy has a set of critical numbers (e.g., one assumption is that penalty costs must be cheaper now than in the future). In both [80] and [83], the classes of demand are essentially the same except for priority. In our case, there may be inherent differences between the classes of demand (e.g., willingness to wait or pay), and we may intentionally backlog customers or reserve inventory for future customers, which further distinguishes how the different classes may be served. In addition, we assume production capacity is limited, we do not make any assumptions on costs over time, and we allow revenue to depend upon customer class.

More recent research in inventory that is relevant includes Sobel and Zhang [73]. In this work, the authors study an inventory problem with fixed plus linear production costs and two demand classes. The deterministic demand class must be satisfied immediately, and the stochastic demand can be backlogged if there is not enough inventory. The main result is that a modified  $(s, S)$  policy is optimal. In our case, our production costs are simpler (linear only), but demand for both classes is stochastic and we allow tactical inventory.

Frank et al. [32] add to the work, again considering one deterministic and one stochastic demand class. They allow the firm to specify how much of the stochastic demand to satisfy; this is somewhat similar to using discretionary sales. Their main result is that a state-dependent optimal policy exists but is quite complex, so they propose a heuristic policy of the form  $(s, k, S)$ , where the rationing policy  $k$  specifies the amount of on-hand inventory to reserve for deterministic demand before ordering; thus,  $k$  also determines the inventory

available to satisfy stochastic demand. Katircioglu and Atkins [47] also consider production and allocation problems with multiple classes of customers. In this work, customer classes require different service levels, and they propose a heuristic that solves the problem myopically and is easy to implement. For our problem, the optimal policy has a simple structure and includes explicit decisions for reserving and backordering (other differences are as outlined above).

One stream of research that considers multiple classes of customers with stochastic demand in manufacturing focuses on *rationing* (see for instance, Dekker et al. [17] or Moon and Kang [61] as well as Topkis [80] reviewed above). The term “rationing” is generally used to refer to the allocation of a resource such as capacity or inventory between competing customer classes. The results in this research area often describe threshold or critical levels that indicate the resource to be allocated to each class. This critical-level policy is optimal for some cases and is used as a heuristic in others. These papers generally focus on dynamic control of a single machine, and they do not consider production problems that span a number of periods with non-stationary parameters. In our case we find threshold values of this type (see the nesting policy for PDS), and we also incorporate resource allocations across time periods.

In most of the described results in the rationing area, a key assumption is that demand is Poisson (see for example, Balakrishnan et al. [5] and Melchioris et al. [57]). In some, there is also an assumption that the production time is exponential (Ha [39]). The most relevant work in this stream is Ha [40], who assumes demand is Poisson and the processing time is Erlang. The key contribution is that the optimal policy has critical levels with monotonic properties. This policy is most similar to the one we find for PDS in this chapter, although in our case we have limited production capacity and tactical inventory. We also consider leadtime differences explicitly and allow planned backlogs.

An important paper that allows tactical inventory is Scarf [68], who introduced discretionary sales into a problem with fixed production setup costs and one customer class. In his case, a base-stock type of policy is optimal for production, but unlike the production decision, the optimal discretionary sales decision should be decided after demand is revealed

in a given period in order to achieve the maximum profit. The use of discretionary sales is also analyzed in Chan et al. [12], which considers a single-class stochastic inventory model with multi-period pricing and production decisions under limited capacity when demand is a general stochastic function.

We build on our work in [12], where we found that a modified base-stock policy with a production and reserving decision pair was optimal, in which the optimal values do not depend on the demand that arrives if price is decided in advance. A fundamental difference in the current research is that we add multiple classes of customers who differ in their willingness to wait (and pay), and we allow delayed fulfillment. The current work also builds on Liu and Simchi-Levi [55], who extended [12] to allow delayed fulfillment until the end of the horizon.

The rest of this chapter is organized as follows. In Section 2.2, we introduce and analyze the Priority Differentiation and Non-Differentiation strategies. We perform computational analysis to compare expected profits under the two strategies in Section 2.3 to explore the effectiveness of market segmentation in manufacturing. Conclusions are contained in Section 2.4.

## ***2.2 Models and Results***

We focus on a single product sold at a single manufacturer over a multi-period time horizon, where the manufacturer has limited production capacity in each period. The manufacturer serves two customer classes, whose demand is ordered by class (i.e., sorted by priority). This means that in any period, first-class demand is fully known by the manufacturer before he has to make a decision regarding second-class demand. The customers of these two classes differ in their priority level and willingness to pay. The first-class customers are willing to pay a premium over the price of the second-class customers in order to have priority access in the current period to both on-hand inventory and backlogging availability. Thus, by paying the premium, first-class customers are satisfied first with the inventory and backlogging resources available by the manufacturer in the current period, and the demand of the second-class customers is addressed with the remaining resources.

The main model that we will consider throughout this chapter is the Priority Differentiation Strategy (PDS), where we assume that the manufacturer has the ability to differentiate the customer classes. We seek to optimize the allocation of limited inventory and production capacity, considering the possibility of reserving inventory to satisfy future demand and allocating future production capacity by backlogging current demand. We show that there is an optimal set of production, backlog, and reserve inventory decisions that allocates current and future resources between customer classes. Considering the general model (PDS) also allows us to analyze other models; for instance, we consider one in which the manufacturer cannot differentiate the customer classes and treats every customer equally (see the Non-Differentiation Strategy (NDS)). This extension and others are described in Section 2.2.3.

### 2.2.1 Notation and Assumptions

The manufacturer makes decisions over a multi-period time horizon,  $t = 1, 2, \dots, T$ , with  $T$  representing the end of the horizon. The production in each period  $t$  is limited by the capacity,  $q_t$ , and the manufacturer pays a production cost per unit of  $c_t$ . Inventory holding cost is linear, and a charge per unit,  $h_t$ , is assessed to carry inventory from  $t$  to  $t + 1$ . Throughout the chapter, the superscripts of 1 and 2 will be used for the first and second classes, respectively.

The manufacturer has predetermined prices,  $p_t^1$  and  $p_t^2$ , for the customers of the first and second classes, respectively, that may be different in each period. Separation of pricing and production decisions is very common in current practice. In some companies, pricing decisions are made by the marketing department before the start of a selling season, while production decisions are made by the operations department.

We assume that each first-class customer is charged a higher price than a second-class customer in the same period; that is,  $p_t^1 > p_t^2$  for each  $t$ , although we make no restrictions on prices between different time periods. This even allows  $p_t^1 < p_{t+1}^2$ , in case there is a significant change in demand curves over time. The salvage value of any units left at the end of the horizon is  $v$ , and  $p_T^1 > p_T^2 > v$ . For classes  $i = 1, 2$ , the cost per unit for demand

in class  $i$  that is rejected and lost is  $\ell_t^i$ , and  $\beta_t^i$  is the cost per unit for demand in class  $i$  that is backlogged. We assume that  $\ell_t^1 > \ell_t^2$  and  $\beta_t^1 > \beta_t^2$  for each  $t$ , since losing or delaying the fulfillment of the first-class customers is more costly than for the second-class customers. We define net revenue of selling to customer class  $i$  from current inventory as  $p_t^i + h_t + \ell_t^i$ ; similarly, the net revenue from backlogging is  $p_t^i - \beta_t^i + \ell_t^i$ . Holding cost,  $h_t$ , is assumed to satisfy  $p_t^1 - \beta_t^1 + \ell_t^1 > p_t^2 + h_t + \ell_t^2$  in each period  $t$ , which ensures that backlogging one first-class customer is more expensive than rejecting a second-class customer to save a unit of inventory for the future.

Each customer belongs to only one demand class, and demand from one class is assumed to be independent of the other class. Each demand function is a general non-stationary stochastic function,  $D_t^i$ , with known probability and cumulative distribution functions  $\phi_t^i$  and  $\Phi_t^i$ , respectively. We assume that the demand function in each period is continuous and differentiable, but no other assumptions are made on the shape of the demand function, so a wide variety of demand models could be used.

Production is a decision made at the beginning of each period and the production leadtime is zero. The net inventory (*on-hand* – *backlogs*) at the beginning of period  $t$  is  $I_t$ , and let  $S_t$  represent the net inventory plus production in period  $t$ . In our initial analysis we restrict ourselves to delivering backordered items one period later, and we assume previously-accepted orders are fulfilled before new orders are accepted, which is possible since we restrict backorders in each period to be no more than the capacity in the next period.

The sequence of events in every period is as follows. At the beginning of a period, the manufacturer checks the inventory level  $I_t$  and decides the production quantity; products arrive immediately, and the manufacturer fulfills the backorders carried from the previous period with the available inventory. Then the demand in the current period is revealed and the manufacturer decides the amount to reserve,  $R_t^i$ , and the amount of future capacity to make available to current customers (i.e., the amount to backlog),  $B_t^i$ .  $R_t^1$  is the amount of inventory to protect from (not sell to) classes 1 and 2, and  $R_t^2$  is the additional amount of inventory to protect from class 2; thus, the total amount to protect from class 2 is  $R_t^1 + R_t^2$ .

**Table 1:** Notation

$q_t$	production capacity in period $t$
$c_t$	production cost per unit in period $t$
$h_t$	inventory holding cost per unit from period $t$ to $t + 1$
$p_t^1$	price charged to first-class customers in period $t$
$p_t^2$	price charged to second-class customers in period $t$
$v$	salvage value of any item left at the end of horizon
$\ell_t^i$	cost per unit for demand in class $i$ that is not satisfied
$\beta_t^i$	cost per unit for demand in class $i$ that is backlogged
$D_t^i$	demand realization of class $i$ in period $t$
$I_t$	net inventory at the beginning of period $t$
$S_t$	net inventory plus production in period $t$
$R_t^1$	amount of inventory to protect from classes 1 and 2
$R_t^2$	amount of additional inventory to protect from class 2
$B_t^2$	amount of future capacity made available to classes 1 and 2
$B_t^1$	amount of additional future capacity made available to class 1

The amount of future capacity to make available to classes 1 and 2 now is  $B_t^2$ , and  $B_t^1$  is the additional capacity for class 1; thus, the total capacity for backlogging class 1 is  $B_t^1 + B_t^2$ . The demand is satisfied according to the  $S_t$ ,  $B_t^i$  and  $R_t^i$  values. The notation that we defined in this section is provided in Table 1 for ease of reference.

### 2.2.2 Priority Differentiation Strategy (PDS)

In the Priority Differentiation Strategy, we assume that the first class is willing to pay a premium to receive priority over all available inventory and backlogging in the current period. The result is that the first and second classes may be fulfilled now or in the next period, depending on the status of the system. Thus, the manufacturer has increased flexibility to match supply and demand.

For the purpose of clarity, we introduce some additional notation in Table 4. Due to our assumption of the ordering of demand classes, we satisfy the first-class demand before the second-class demand. Consequently the available inventory for the second class is limited by the first-class demand that is realized. We define the amount of inventory available after the first-class demand is satisfied as  $S_t^2$ . Since the first class has higher priority in the current period, we use as much of  $B_t^2$  as necessary to backlog the first-class demand. Then we use the remaining part of  $B_t^2$  (if there is any left) to backlog the second-class

**Table 2:** Additional notation

$S_t^2$	$= (S_t - R_t^1 - D_t^1)^+$	available inventory after first-class
$B_t^{2,ef}$	$= [B_t^2 - [D_t^1 - [S_t - R_t^1]^+]^+]^+$	<i>effective</i> backlog amount after first-class
$A_t^1$	$= \min(B_t^1 + B_t^2, [D_t^1 - [S_t - R_t^1]^+]^+)$	actual backlogged orders from first class
$A_t^2$	$= \min(B_t^{2,ef}, [D_t^2 - S_t^2 + R_t^2]^+)$	actual backlogged orders from second class
$I_{t+1}^{R^1}$	$= \min(S_t, R_t^1)$	inventory carried forward due to $R^1$ decision
$I_{t+1}^{R^2}$	$= \min(S_t^2, R_t^2)$	inventory carried forward due to $R^2$ decision
$I_{t+1}^{low}$	$= [S_t^2 - R_t^2 - D_t^2]^+$	inventory carried forward due to low demand

demand. We call this remaining backlog availability  $B_t^{2,ef}$ , or the *effective* backlog amount after first-class demand. For ease of presentation in the chapter, we further define the actual backlogged orders from first and second-class customers after demand is satisfied as  $A_t^1$  and  $A_t^2$ , respectively, and inventory carried forward due to  $R^1$  and  $R^2$  as  $I_{t+1}^{R^1}$  and  $I_{t+1}^{R^2}$ , respectively. If demand is *low* enough so that there is leftover inventory at the end of the period, we denote this additional inventory as  $I_{t+1}^{low}$  (see Table 2 for summary of the additional notation).

We model the PDS problem as a Markov decision process, where the state of the system is represented by the net inventory. For clarity of exposition, we present the model with the  $R_t^i$  and  $B_t^i$  decisions given *ex ante*. However, in our analysis we show that the optimal  $R_t^{i*}$  and  $B_t^{i*}$  decisions are the same whether they are made before or after demand revelation. Let  $J_t(I_t)$  be the expected profit from period  $t$  forward to the end of the horizon, or the *profit-to-go*. Let  $G_t(S_t)$  be the expected profit-to-go with  $S_t$  units of product available after production. The first and second derivatives of  $J_t(I_t)$  are denoted, respectively, as:  $J_t'(I_t)$  and  $J_t''(I_t)$ ; the derivatives of other functions are indicated similarly. We can now write the optimal expected profit in period  $t$  and onward for the PDS problem as the following recursive equation.

$$J_t(I_t) = \max_{S_t: \max(0, I_t) \leq S_t \leq I_t + q_t} \{-c_t(S_t - I_t) + G_t(S_t)\}, \quad \text{where} \quad (1)$$

$$\begin{aligned}
G_t(S_t) = \max_{B_t^1, B_t^2, R_t^1, R_t^2} \int \int \bigg\{ & p_t^1 \min(D_t^1, S_t - R_t^1 + B_t^1 + B_t^2) \\
& + p_t^2 \min(D_t^2, [S_t^2 - R_t^2]^+ + B_t^{2,ef}) \\
& - h_t [S_t^2 - R_t^2 - D_t^2]^+ - h_t \min(R_t^2, S_t^2) - h_t R_t^1 \\
& - \ell_t^1 (D_t^1 - [S_t - R_t^1]^+ - B_t^1 - B_t^2)^+ \\
& - \ell_t^2 (D_t^2 - [S_t^2 - R_t^2]^+ - B_t^{2,ef})^+ \\
& - \beta_t^1 \min([D_t^1 - S_t + R_t^1]^+, B_t^1 + B_t^2) \\
& - \beta_t^2 \min([D_t^2 - [S_t^2 - R_t^2]^+]^+, B_t^{2,ef}) \\
& + J_{t+1}((I_{t+1}^{low} + I_{t+1}^{R^1} + I_{t+1}^{R^2}) - A_t^1 - A_t^2) \bigg\} d\Phi_t^1(D_t^1) d\Phi_t^2(D_t^2), \quad (2)
\end{aligned}$$

$$\text{subject to: } B_t^1 + B_t^2 \leq q_{t+1}, \quad R_t^1 + R_t^2 \leq S_t.$$

In Equation (1), the maximization of profit is over the target inventory decision. The first term of the function is the production cost; the production also covers any backlogged orders from the prior period. The second term is the profit in the remainder of the period (and horizon) starting with the available inventory after production is completed and backorders are fulfilled.

In Equation (2), the function  $G_t$ , the profit-to-go after production, is maximized over the reserve inventory and backlogging decisions. The first element of the function is the revenue from first-class customers, including both physical inventory and backlogged orders. The second term is revenue from the second-class demand with available inventory and backlogged orders. The third piece is the inventory holding cost to be paid for all inventory not sold. The fourth and fifth terms represent inventory holding cost that is incurred for all inventory reserved for the future. The sixth and seventh terms are the rejection penalties for demand not satisfied for the first and second classes, respectively, and the eighth and ninth terms are the delay penalty associated with the backlogged demand for the first and second classes, respectively. The last term in the equation represents the profit in future periods, sending forward any leftover physical inventory and backlogged orders. For period  $T$ , the final term is replaced by the salvage cost of leftover inventory, namely  $v(S_T - D_T^1 - D_T^2)^+$ . Finally, the constraints ensure that the manufacturer does not sell more future capacity



than he has or reserve more inventory than is available.

### 2.2.2.1 Problem Simplifications

For each demand class, the manufacturer decides the amount of inventory to reserve and the amount of backordering. To simplify the problem at hand, we show that in an optimal policy for a class and a time period, at least one set of these decisions must be zero.

**Lemma 2.1.** *In any optimal policy under the Priority Differentiation Strategy, we have:*

$$(B_t^1 + B_t^2) \cdot R_t^1 = 0 \text{ and } B_t^2 \cdot (R_t^1 + R_t^2) = 0 \quad t = 1, 2, \dots, T.$$

The first of these conditions says that if it is good to protect items for the future from class 1 and lose some of the current demand, then it is not reasonable to backorder items from class 1 or the lower-revenue class 2 (the contrapositive is also true). Likewise, the second condition says that if it is good to backorder demand from even the (lower-paying) second class in the current period, then it will not be reasonable to protect items from (and lose demand from) the second class or the higher-paying first class in the current period (the contrapositive is also true). The formal proof can be found in the Appendix.

By Lemma 2.1, the structure of the optimal policies can be simplified. In each period there are three candidate policies, of which the best policy will be chosen; this choice will be dependent on the state of the system. The possible options are to Reserve-Inventory ( $R_t^1 \geq 0, R_t^2 \geq 0$ ), to Backlog-Demand ( $B_t^1 \geq 0, B_t^2 \geq 0$ ), or to Reserve-and-Backlog ( $R_t^2 \geq 0, B_t^1 \geq 0$ ). Thus,

$$G_t(S_t) = \max\{G_t^1(S_t), G_t^2(S_t), G_t^3(S_t)\},$$

where  $G_t^1(S_t)$ ,  $G_t^2(S_t)$ , and  $G_t^3(S_t)$  represent the profit-to-go with  $S_t$  units of products available after production under the Reserve-Inventory policy, the Backlog-Demand policy and the Reserve-and-Backlog policy, respectively. These three policies are given by:

$$\begin{aligned} G_t^1(S_t) &= \max_{R_t^1 + R_t^2 \leq S_t} \left\{ g_t^1(S_t, R_t^1, R_t^2) \right\}, \\ G_t^2(S_t) &= \max_{B_t^1 + B_t^2 \leq q_{t+1}} \left\{ g_t^2(S_t, B_t^1, B_t^2) \right\}, \\ G_t^3(S_t) &= \max_{B_t^1 \leq q_{t+1}, R_t^2 \leq S_t} \left\{ g_t^3(S_t, R_t^2, B_t^1) \right\}. \end{aligned}$$

In each of the three cases, the starting inventory after production is completed and backorders are fulfilled is  $S_t$ . The first function,  $g_t^1(S_t, R_t^1, R_t^2)$ , indicates the profit-to-go when inventory may be protected from both classes ( $R_t^1, R_t^2 \geq 0$ ). In this case the manufacturer will not backlog orders of current customers because the backlog orders will reduce the future capacity available to customers (therefore  $B_t^1 = B_t^2 = 0$ ). The profit from this policy is represented as:

$$g_t^1(S_t, R_t^1, R_t^2) = \iint \left\{ p_t^1 \min(D_t^1, S_t - R_t^1) - \ell_t^1 (D_t^1 - S_t + R_t^1)^+ + p_t^2 \min(D_t^2, [S_t^2 - R_t^2]^+) \right. \\ \left. - \ell_t^2 (D_t^2 - [S_t^2 - R_t^2]^+)^+ - h_t [S_t^2 - R_t^2 - D_t^2]^+ - h_t \min(R_t^2, S_t^2) \right. \\ \left. - h_t R_t^1 + J_{t+1} (I_{t+1}^{low} + I_{t+1}^{R^1} + I_{t+1}^{R^2}) \right\} d\Phi_t^1(D_t^1) d\Phi_t^2(D_t^2).$$

The function  $g_t^2(S_t, B_t^1, B_t^2)$  indicates the profit-to-go when backorders for each class may be desirable ( $B_t^1, B_t^2 \geq 0$ ). However, the manufacturer will not protect inventory from either class ( $R_t^1 = R_t^2 = 0$ ). The resulting formulation is:

$$g_t^2(S_t, B_t^1, B_t^2) = \iint \left\{ p_t^1 \min(D_t^1, S_t + B_t^1 + B_t^2) + p_t^2 \min(D_t^2, [S_t - D_t^1]^+ + B_t^{2,ef}) \right. \\ \left. - \ell_t^1 (D_t^1 - S_t - B_t^1 - B_t^2)^+ - \ell_t^2 (D_t^2 - [S_t - D_t^1]^+ - B_t^{2,ef})^+ \right. \\ \left. - \beta_t^1 \min([D_t^1 - S_t]^+, B_t^1 + B_t^2) - \beta_t^2 \min([D_t^2 - [S_t - D_t^1]^+]^+, B_t^{2,ef}) \right. \\ \left. - h_t [S_t - D_t^1 - D_t^2]^+ + J_{t+1} (I_{t+1}^{low} - A_t^1 - A_t^2) \right\} d\Phi_t^1(D_t^1) d\Phi_t^2(D_t^2).$$

The remaining function,  $g_t^3(S_t, R_t^2, B_t^1)$ , indicates the profit-to-go when the manufacturer may backlog orders of the first class for future fulfillment ( $B_t^1 \geq 0$ ) and may also protect inventory from the second class for future use ( $R_t^2 \geq 0$ ).

$$g_t^3(S_t, R_t^2, B_t^1) = \iint \left\{ p_t^1 \min(D_t^1, S_t + B_t^1) + p_t^2 \min(D_t^2, [S_t^2 - R_t^2]^+) - h_t [S_t^2 - R_t^2 - D_t^2]^+ \right. \\ \left. - \beta_t^1 \min([D_t^1 - S_t]^+, B_t^1) - \ell_t^1 (D_t^1 - S_t - B_t^1)^+ - \ell_t^2 (D_t^2 - [S_t^2 - R_t^2]^+)^+ \right. \\ \left. - h_t \min(R_t^2, S_t^2) + J_{t+1} ((I_{t+1}^{low} + I_{t+1}^{R^2}) - A_t^1) \right\} d\Phi_t^1(D_t^1) d\Phi_t^2(D_t^2).$$

In each period one of these three policies will be chosen, and this choice also impacts the future state of the system. Intuition gives us some idea of when each policy will be selected, which we establish more formally in our results below. We expect that the Reserve-Inventory policy will be selected in a period where the marginal expected profit from selling each of

the reserve units in the future is better than the net revenue of selling a unit now out of inventory. For the Backlog-Demand policy, intuition suggests that it will be best when the net revenue of backlogging in the current period is better than the marginal expected profit from selling each of the units in the future. Finally, the Reserve-and-Backlog policy will be optimal when the net revenue of backlogging to the first-class customers is significantly greater than the marginal expected future profit of the backlogged units, but the second class has a lower net revenue when selling from inventory than the marginal expected future profit of sending forward reserved units.

### 2.2.2.2 Results

Under the Priority Differentiation Strategy, we can show that all the profit-to-go functions have nice structure (quasi-concave or concave), thus yielding easy to implement decisions. These results are summarized in the following theorem (see the Appendix for the full details):

**Theorem 2.1.** *Under the Priority Differentiation Strategy,*

- $g_t^1(S_t, R_t^1, R_t^2)$  is a quasi-concave function of  $R_t^1$  and  $R_t^2$ , for all  $t = 1, \dots, T$ .
- $g_t^2(S_t, B_t^1, B_t^2)$  is a quasi-concave function of  $B_t^1$  and  $B_t^2$ , for all  $t = 1, \dots, T$ .
- $g_t^3(S_t, B_t^1, R_t^2)$  is a quasi-concave function of  $B_t^1$  and  $R_t^2$ , for all  $t = 1, \dots, T$ .
- $G_t(S_t)$  is a concave function of  $S_t$ , for all  $t = 1, \dots, T$ .
- $J_t(I_t)$  is a concave function of  $I_t$ , for all  $t = 1, \dots, T$ .
- The unconstrained optimizers  $(R_t^{1*}, R_t^{2*}, B_t^{1*}, \text{ and } B_t^{2*})$  for functions  $g_t^1(S_t, R_t^1, R_t^2)$ ,  $g_t^2(S_t, B_t^1, B_t^2)$  and  $g_t^3(S_t, B_t^1, R_t^2)$ , are independent of inventory level  $S_t$  and demand realizations  $D_t^1$  and  $D_t^2$ , where

$$\begin{aligned} (R_t^{1*}, R_t^{2*}) &\in \arg \max_{0 \leq R_t^1, 0 \leq R_t^2} \left\{ g_t^1(S_t, R_t^1, R_t^2) \right\}, \\ (B_t^{1*}, B_t^{2*}) &\in \arg \max_{0 \leq B_t^1, 0 \leq B_t^2} \left\{ g_t^2(S_t, B_t^1, B_t^2) \right\}, \\ (R_t^{2*}, B_t^{1*}) &\in \arg \max_{0 \leq R_t^2, 0 \leq B_t^1} \left\{ g_t^3(S_t, B_t^1, R_t^2) \right\}. \end{aligned}$$

In Section 2.2.1, while explaining the sequence of events, we assumed that the  $R_t^i$  and  $B_t^i$  decisions are made after seeing the demand. In Theorem 2.1 we show that these decisions are independent of the demand in the current period; thus, the manufacturer can decide their optimal levels before the demand is revealed for the period. The theorem implies the optimal policy for the Priority Differentiation Strategy; thus, we have the following corollary.

**Corollary 2.1.** *Given a vector of prices, there exists an optimal modified base-stock policy for the Priority Differentiation Strategy with an optimal order-up-to level ( $S_t^*$ ), and for  $i = 1, 2$  optimal reserve-up-to-levels ( $R_t^{i*}$ ) and optimal backlog-up-to levels ( $B_t^{i*}$ ).*

We refer to the policy as *modified* base-stock because it may be limited by capacity or available inventory. If there is not sufficient capacity to bring the inventory level up to  $S_t^*$ , then as much as possible should be produced. Similarly, the  $R_t^i$  and  $B_t^i$  decisions are limited by  $S_t$  and  $q_{t+1}$ , respectively. The form of the optimal decisions are apparent from the concavity and quasi-concavity of the profit functions. At each stage in the problem, the manufacturer trades off the current net revenue against the marginal future contribution in terms of cost or revenue and chooses the best allocation of resources.

Additional insight may be gained by looking at the optimal decisions in more detail. The optimal decisions are defined by the following:<sup>1</sup>

$$\begin{aligned}
S_t^* &= \max\{S : c_t \leq G'_t(S)\} && \text{if } c_t \leq G'_t(0) \\
R_t^{1*} &= \max\{I : p_t^1 + \ell_t^1 + h_t \leq J'_{t+1}(I)\} && \text{if } p_t^1 + \ell_t^1 + h_t < J'_{t+1}(0) \\
R_t^{1*} + R_t^{2*} &= \max\{I : p_t^2 + \ell_t^2 + h_t \leq J'_{t+1}(I)\} && \text{if } p_t^2 + \ell_t^2 + h_t < J'_{t+1}(0) \\
B_t^{1*} + B_t^{2*} &= \min\{I : J'_{t+1}(-I) \geq p_t^1 + \ell_t^1 - \beta_t^1\} && \text{if } p_t^1 + \ell_t^1 - \beta_t^1 > J'_{t+1}(0) \\
B_t^{2*} &= \min\{I : J'_{t+1}(-I) \geq p_t^2 + \ell_t^2 - \beta_t^2\} && \text{if } p_t^2 + \ell_t^2 - \beta_t^2 > J'_{t+1}(0).
\end{aligned}$$

In Figure 1 we show the marginal expected profit in period  $t + 1$  as a function of inventory. According to the decisions described above, an optimal decision (e.g., the reservation decision  $R_t^{1*}$ ) equals the inventory level where the relevant prices and costs (e.g.,  $p_t^1 + \ell_t^1 + h_t$ )

---

<sup>1</sup>Each decision is equal to 0 if the condition is never satisfied.

cross the marginal expected profit curve. Figure 1(a) illustrates the reserve inventory decisions, which correspond to the Reserve-Inventory policy in the previous section.

In the remaining figures, we show the marginal expected profit curve compared to the costs relevant to the other optimal decisions above. The optimal backlogging decision is portrayed in Figure 1(b); this decision corresponds to the Backlog-Demand Policy. Finally, we show the optimal decision that results from Reserve-and-Backlog Policy in Figure 1(c). A similar picture could be drawn for the target inventory decision comparing the production cost ( $c_t$ ) with the derivative of the  $G_t$  function; this is left out for brevity. In all of the decisions, we note that the manufacturer is trading off the certain net revenue in the current period (e.g.,  $p_t^1 + \ell_t^1 + h_t$ ) with a marginal expected profit in the future. Clearly there is some risk with betting on the future, but such trade-offs are made regularly in many situations.

### 2.2.3 Special Cases and Extensions

We are also interested in situations in which manufacturers cannot differentiate customers and treat them as a single class. We denote this situation as the Non-Differentiation Strategy (*NDS*), which is a special case of PDS. We assume that the manufacturer takes the second-class customers' reservation price,  $p_t^2$ , as the selling price to all customers. Since the lower price is charged to both classes, customers in both classes are willing to wait one period if the item is not available to them, as in PDS. The difference of NDS from PDS is not in the customers' preferences, but in the manufacturer's treatment of the customers. First-class customers would be willing to pay extra if the manufacturer could differentiate, but he is not able to or willing to differentiate. If we set  $D_t^2$ ,  $R_t^2$ , and  $B_t^2$  to zero and replace  $D_t^1$  with total demand in the formulation of PDS, we get NDS. Thus, the optimal policy is of the form  $(S, R, B)$  as in PDS.

Initially, we analyzed the PDS problem for two customer classes under the assumption that all backlogged orders are filled within one period. However, there are several more general extensions that easily follow from our initial proof. Some of these extensions are outlined below.

- *Multiple classes:* Our results for PDS hold for a problem with more than two customer classes. As before, it is necessary to assume that each class has priority in the allocation of inventory and production capacity over the lower priority classes in the current period. With this assumption, the nesting structure of the tactical inventory decisions is still optimal. To be more specific, one could have a menu (price, priority ranking) for each customer class. If there are many customer classes, it might be difficult for customers to choose from the sets, and at the firm level, priority ordering of many classes would also be difficult. However, it may be reasonable for 3 - 5 classes, which can occur in some applications.
- *Time-differentiated customers:* It is also possible to extend the models to cover situations where some classes are always served immediately while others receive immediate or delayed fulfillment. An example for this is a Time Differentiation Strategy (TDS), where the first-class customers would never be willing to wait and are served immediately, while the second-class customers can be served immediately or next period. For this problem, the optimal policy is in the form of  $(S, R^i, B)$  for  $i = 1, 2$ , which is a critical threshold policy as before. See Chapter 3 for details.
- *Long Leadtime:* The fulfillment leadtime in our analysis is assumed to be one period. However, it is also possible to allow for planned backlogs where the orders can be delivered anytime before the end of the time horizon. For the extended analysis, we assume that backlogs must be filled before new orders are accepted, and under this assumption our nested threshold policies are still optimal. If there is a leadtime  $1 < l_t < T - t$  in each period that specifies orders must be delivered in period  $t + l_t$ , then the problem is structurally more complex.<sup>2</sup> In particular, the state space increases since previous orders must be tracked so that they are fulfilled in the correct time period. Furthermore, even if the expected profit is concave, the optimal policy may be complex and not easy to implement.

---

<sup>2</sup>If  $l$  indicates that orders must be filled *by* period  $t + l$ , and previously-accepted orders must be filled before new orders are accepted, then the results in this paper hold as described for planned backlogs. However, for the version of the problem with specific and varying  $l_t$ , the assumption that previous orders are filled first may be too restrictive.

### 2.3 Computational Analysis

In this section we report on a computational study conducted to obtain insights about the benefits of customer differentiation and tactical inventory use in PDS and NDS. Our goal is to examine the relative performance of the policies of the  $(S, R, B)$  form in different problem settings and identify the situations where this type of policy can provide significant increases in profit.

The benchmark we use for each of our strategies is a traditional base-stock policy where the manufacturer uses the modified order-up-to policy ( $S$  policy) and serves all customers as in a single class. We assume that sales are lost if there is insufficient inventory on-hand or if customers are rejected. We compare the performance of the  $(S, R, B)$  type policies over the traditional policies using the metric of profit potential, as defined by  $100 * (\frac{V_{(S,R,B)}}{V_S} - 1)$ , where  $V$  indicates the expected profit of the problem being solved. In both the traditional policy and NDS, we use  $p_t^2$  as the price charged to all customers to ensure that we serve to both of the classes. This implies that PDS may show a big improvement in profit that is due, in part, to the ability to differentiate customers.

The profit improvement of the Priority Differentiation Strategy compared to traditional inventory policies comes from three sources: prioritized demand classes, differentiated pricing, and shifting inventory to the next period, whereas the Non-Differentiation Strategy only has the last source. Thus, by comparing both PDS and NDS to the traditional policy, we can separate the impact of price differentiation versus tactical inventory.

#### 2.3.1 Experiment Details

The total average demand from the first and second-class customers equals 100 in each experiment. We assume that demand uncertainty is additive with a mean of 0. We define the coefficient of variation of demand in a given period as  $CV_U^i = s(D_t^i)/E(D_t^i)$ , where  $s$  denotes the standard deviation, and  $E$  denotes the expected value. In all cases shown, the coefficient of variation of demand uncertainty is the same in each period and is equal to 0.2.

Production capacity is constant for a particular instance, while it is allowed to take the values of 60% (low), 80% (med), and 100% (high) of the expected total average demand

**Table 3:** Specific experimental data

<b>t</b>	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>	<b>5</b>	<b>6</b>	<b>7</b>	<b>8</b>	<b>9</b>	<b>10</b>	<b>11</b>	<b>12</b>	<b>Avg</b>
$c_t$	70	90	70	50	70	90	70	50	70	90	70	50	<b>70</b>
$p_t^2$	90	110	90	70	90	110	90	70	90	110	90	70	<b>90</b>
$p_t^1$	110	130	110	90	110	130	110	90	110	130	110	90	<b>110</b>

for both classes over the horizon (denoted by  $Dem^*$ ) in some experiments. The production cost may vary by period, but the production cost vector is the same across instances. (We also ran experiments where the production cost is the same in each period and obtained similar results.) See Table 3 for the exact data; for example, the average markup of  $p_t^2$  ( $p_t^1$ ) over the cost is about 30%(60%) for the experiments on class proportions.

We study the impact of the percentage of first versus second-class demand in our first set of experiments. In these cases the expected demand from first-class customers over the horizon,  $E(D^1)$ , takes the values of 20, 25, 50, 75, and 80, and the expected second-class demand,  $E(D^2)$ , equals  $100 - E(D^1)$ . The prices are constant over the set of experiments but may vary by period. Having varying prices increases the likelihood that all of the policies will be optimal in some period of an experiment, since the prices create an incentive to shift capacity. The average ratio of  $p_t^1/p_t^2$  is 1.22 for the experiments studying the proportion of demand. See Table 3 for the prices used in this set of experiments.

We also consider the relative price difference between classes. In these experiments  $E(p^2)$  is fixed over the instances, and the price for the first class is set according to  $E(p^1)/E(p^2) = 1.1, 1.2, \text{ and } 1.3$ , where  $E(p^i)$  represents the average price over the horizon. We allow the trend of  $p_t^2$  (and correspondingly,  $p_t^1$ ) to be either linearly increasing or decreasing (we also ran experiments with no clear price trend). Let  $\gamma = p_{t+1}^2 - p_t^2$ , which we assume to be fixed for all  $t = 1 \dots T - 1$ ;  $\gamma$  shows the rate of change of price over time. For the increasing price experiments,  $p_1^2 = \$70$ , and for the decreasing price experiments,  $p_{12}^2 = \$70$  where 12 is the last period.



### 2.3.2 Results

For all experiments, the policies in PDS and NDS using tactical inventory have a higher profit than the traditional policy. This is clear because of the usage of  $p_t^2$  for all customers in the traditional policy. However, note that the profit difference is significant, even when the  $E(D^1)$  percentage is small (see Figure 2(a)).

The performance for a given proportion of first-class customers is better under the tactical inventory policies when the capacities are tight. As an example, in Figure 2(a) the performance of PDS when capacity is  $0.6 Dem^*$  is better than the performance of PDS when capacity is  $0.8 Dem^*$ . As expected for a given capacity level, the performance of the tactical inventory policy in PDS increases almost linearly as the proportion of first-class customers increases. This profit improvement is due to the additional revenue opportunities that the tactical inventory policies have over the traditional policy including higher revenue from first-class customers and an increased ability to meet demand by shifting capacity.

As expected, the profit under NDS is insensitive to the first-class proportion since it does not differentiate between the classes. However, the significant profit improvement over the traditional policy, even though both NDS and the traditional policy offer  $p_t^2$  to everyone, suggesting that the tactical inventory may greatly improve profit. In our experiments, production cost and prices are time varying and capacity is limited. When all parameters are stationary over time and there is sufficient production capacity, the differentiation strategies are unlikely to offer as much improvement over the traditional policy.

For several levels of price proportions ( $E(p^1)/E(p^2)$ ), we look at the rate of price increase (measured by  $\gamma$ ) over the time horizon in Figure 2(b); the decreasing price trend showed nearly the same results. Whether or not the pricing trend is increasing or decreasing, the performance of the  $(S, R^i, B^i)$  policy relative to the traditional policy increases with decreasing  $\gamma$ . To see this for the case of increasing prices, note in Figure 2(b) that the performance of PDS when  $\gamma = 1$  is better than the performance of PDS when  $\gamma = 4$  at every ratio of price differences between the classes. This result is somewhat surprising. Looking at our results more closely, we find that as  $\gamma$  increases, the profits of PDS and the traditional policy are both increasing because the mean prices are increasing. In fact,

the absolute profit difference between the two strategies is increasing with  $\gamma$ . However, the *percentage* profit improvement is not increasing with  $\gamma$ . This seems to be because the total demand in each case is constant and the additional marginal profit from selling one more unit in PDS is small relative to the overall increase in the profit of the traditional policy when  $\gamma$  is large.

The values of the average tactical inventory levels for PDS and NDS are depicted in Figures 3(a) and 3(b), respectively, for increasing prices. In the increasing price experiments all three policies in PDS are active, while for decreasing prices (not shown) only the Backlog-Demand Policy resulted. In some cases the magnitude of the average tactical inventory increases with  $\gamma$  (that is, with increasing trend in price), but this is not true in all cases. Note here that Figure 3 depicts average tactical inventory over the horizon, not necessarily in each period. When we look at the solutions in more detail for increasing price trend, we find that the reserve inventory is used in periods with lower prices and backlogging is used in periods with higher prices. Thus, for each  $\gamma$  level in Figure 3, we have positive backlogging. This is also due to the fact that the backlogging decision is comparing the net revenue from a certain current customer with the marginal *expected* profit from a future customer. In NDS all available tactical inventory decisions are employed, and in some cases (e.g.  $\gamma = 6$ ), the best value of  $R$  for NDS is approximately equal to  $R_t^1 + R_t^2$  in PDS, suggesting that NDS is partially compensating for limited flexibility with high values of tactical inventory for the single customer class.

In the experiments thus far, we set the regular price to be  $p_t^2$  (in the traditional policy), and some customers are willing to pay a higher price  $p_t^1$  for priority service (in PDS) over no priority at regular price  $p_t^2$ . In this situation PDS clearly offers an advantage over the traditional policy, since the average revenue is higher. However, it is also interesting to see what happens when the regular price is  $p_t^1$  and some customers are willing to be served at a lower priority for a discount, paying  $p_t^2$ . We show the results of these experiments with increasing first-class demand in Figure 4(a), where the total number of expected customers is 100 as before. Note here that PDS does not necessarily provide an improvement over the traditional policy, since the average revenue per customer is less than in the traditional

policy. When the first-class demand proportion is more than 50% and capacity is tight, we see that PDS can have a higher profit than the traditional policy, even though the latter has larger average revenue. This result suggests that tactical inventory to shift capacity can overcome the average revenue decrease per customer in some cases. We also consider experiments where the regular price (in the traditional policy) is the average of  $p_t^1$  and  $p_t^2$  in each period, and PDS has some customers willing to pay more ( $p_t^1$ ) for higher priority and some customers willing to have a lower priority for a price discount ( $p_t^2$ ). In this case we see that PDS has greater profit than the traditional model in almost all cases, even though the average revenue is less in PDS when the expected first-class demand is less than 50% (see Figure 4(b)). The main insight from these graphs is that if prioritization of demand classes costs a firm in the average revenue per customer, the benefit of tactical inventory may outweigh the revenue loss.

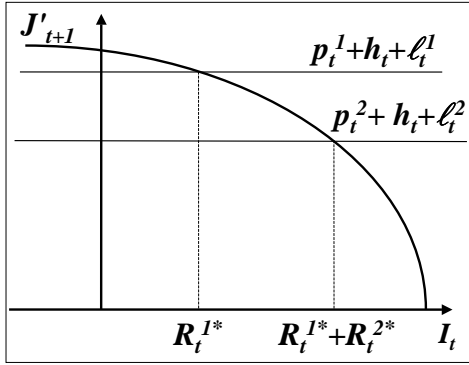
## 2.4 Conclusions

In this chapter we analyzed a multiple-class customer problem where production and tactical inventory decisions must be made in every period and demand is a general stochastic function of time and customer class. We have shown that there are a variety of problems using tactical inventory decisions for which a threshold policy in each period is optimal under a Priority Differentiation Strategy. Specifically, we have a modified base-stock policy consisting of the target inventory decision ( $S$ ), the reserve-up-to levels ( $R^i$ ), and the backlog-up-to levels ( $B^i$ ) for each demand class, or an  $(S, R^i, B^i)$  policy. Under prioritized demand this policy is further nested by customer class.

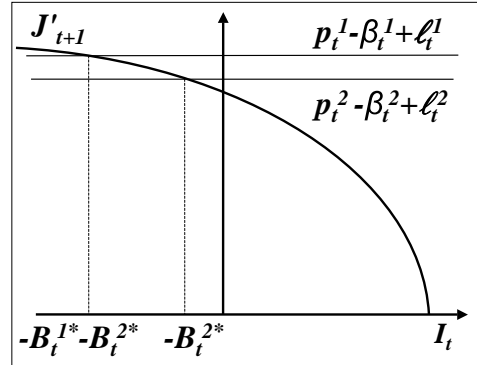
The problem we model and analyze may also have application in other industries. For instance, in some healthcare environments there may be multiple customer classes competing for time on a piece of equipment where priorities are based on the status of the illness. In this problem, there may not be an explicit production decision, but one could still apply backordering and reservation decisions such as promising to service a lower priority class customer in a future time period.

Clearly the analysis in this chapter makes assumptions to simplify the problem, such

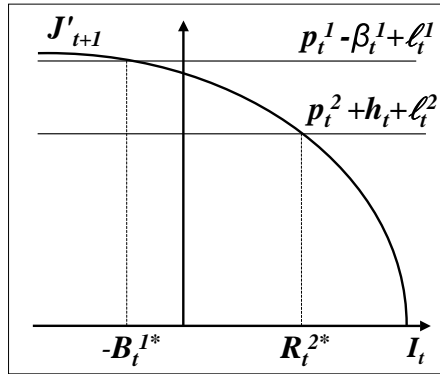
as focusing on a single product. Yet these simplifications allow the development of an optimal policy that is easy to understand, and more importantly, is easy to implement; and the results have extensions beyond those focused on in this article. Further, the simple structure of the threshold policy may give insight for policies to apply to more complicated problems.



(a) Reserve-Inventory Policy

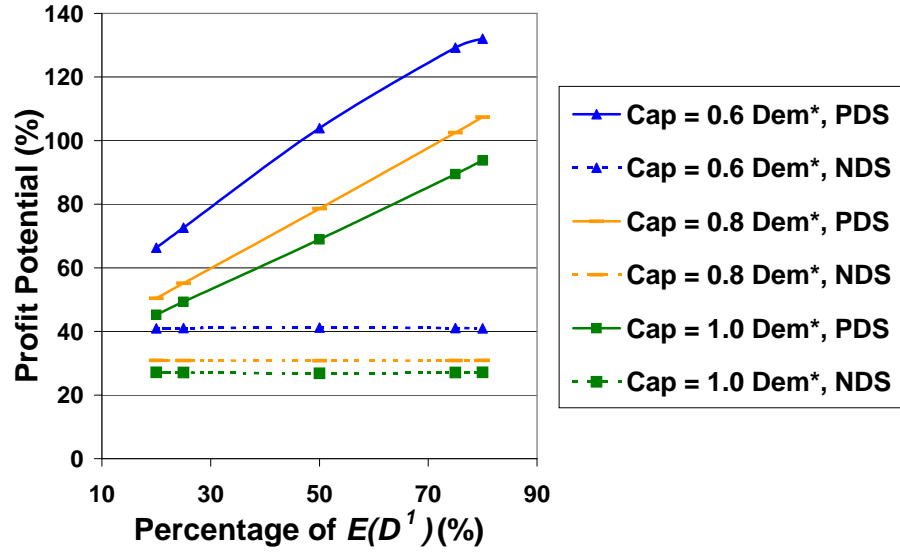


(b) Backlog-Demand Policy

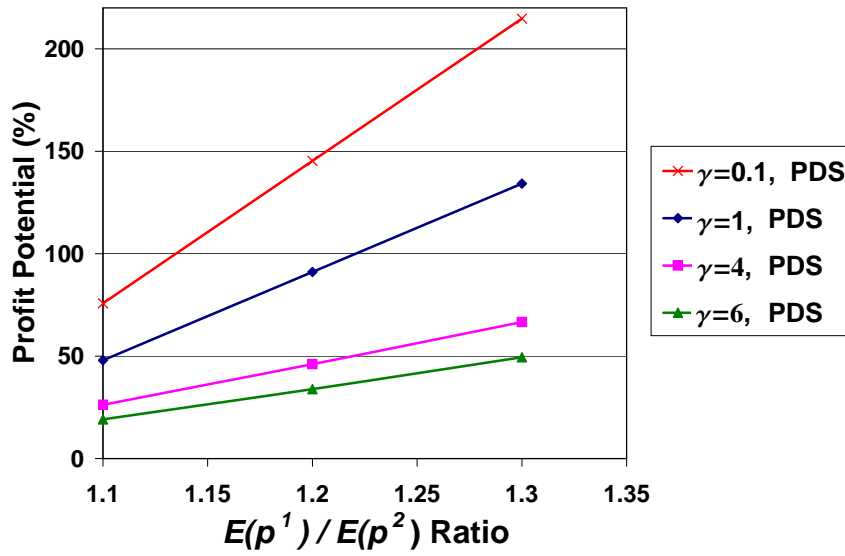


(c) Reserve-and-Backlog Policy

**Figure 1:** Optimal Decisions under the Optimal Policies for PDS

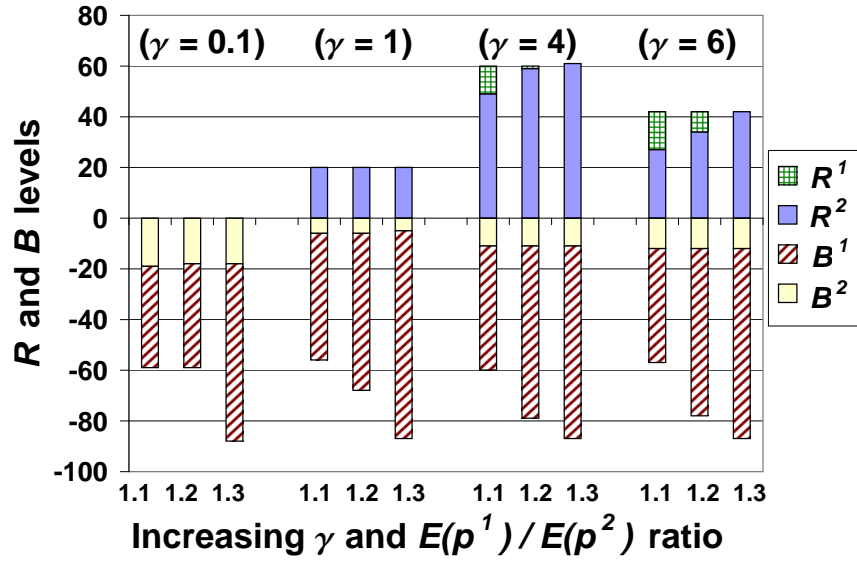


(a) Impact of first-class proportion

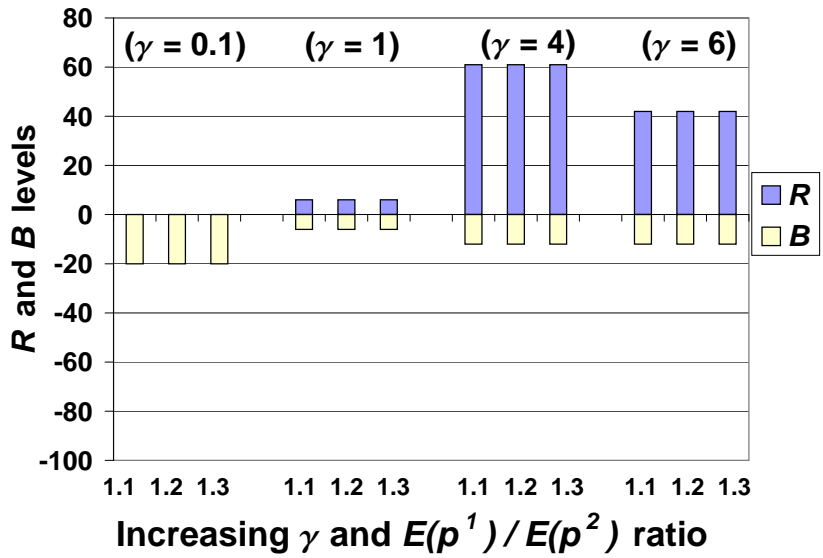


(b) Impact of rate of price increase

**Figure 2:** The relative performance of ( $S$ ,  $R$ ,  $B$ ) type policies over traditional

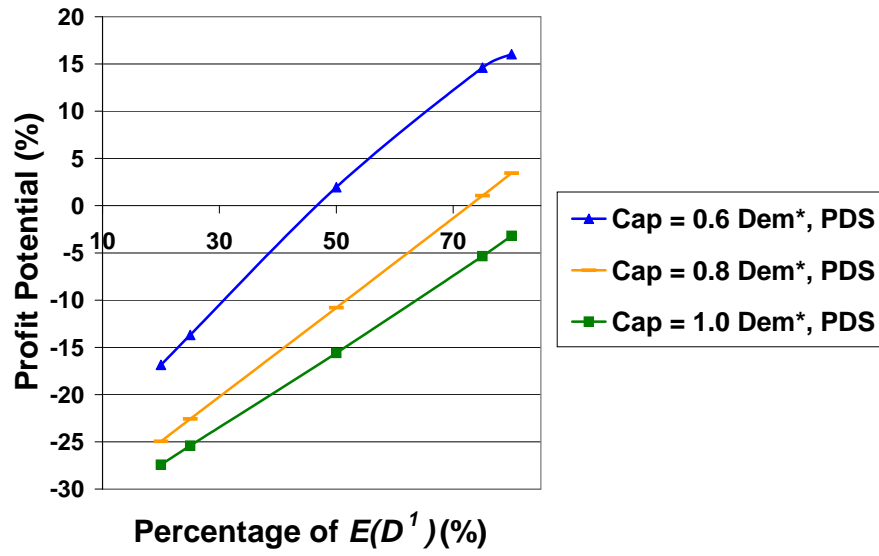


(a) PDS

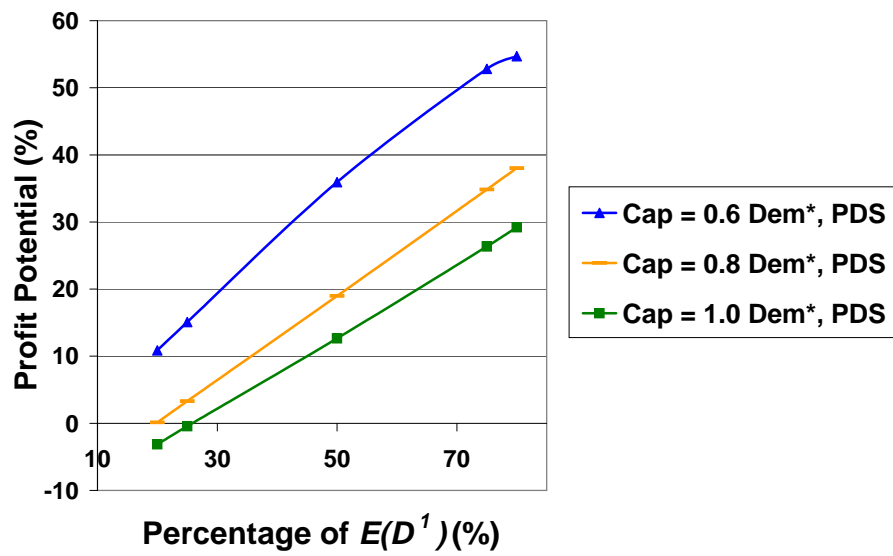


(b) NDS

Figure 3: Tactical Inventory Levels with Price Increasing over the Horizon



(a) Price for traditional is  $p_t^1$



(b) Price for traditional is average of  $p_t^1$  and  $p_t^2$

Figure 4: The relative performance of PDS over traditional with different prices



## CHAPTER III

### POLICIES UTILIZING TACTICAL INVENTORY FOR SERVICE DIFFERENTIATED CUSTOMERS

#### *3.1 Introduction*

Some manufacturing or retail companies now segment their customers according to service and price, since it is not uncommon for customers to have different service and price preferences and the differentiation may benefit the firms. For example, one class of customers may be given immediate fulfillment while another class might receive delayed fulfillment for a discount. For instance, if an executive's laptop has been stolen he may pay a premium for immediate delivery, while someone ordering a computer to go to college may order in advance for a discount. Amazon.com also offers price and delivery time options where paying a price premium gives a customer immediate fulfillment while receiving a Super Saver Shipping discount gives Amazon the opportunity for delayed fulfillment. This may provide greater customer utility (either increased service or decreased price as desired by different customers), while offering greater flexibility to the firm in managing the production system.

Though this can increase utility to the customer or the firm, it is necessary to analyze how to manage the system, which may be more complicated due to the service differentiation. One method to manage this kind of system is to use *tactical inventory*, where current inventory may be set aside to satisfy future demand, and delayed fulfillment of current customers (or "backlogs") may be planned. Tactical inventory may increase profits while ensuring that service of both kinds of customers is met.

The use of tactical inventory is considered in Scarf [68], in which the idea of protecting inventory from being sold to current customers or "discretionary sales" is introduced. Scarf showed that a base-stock policy is optimal for a single-class problem when production setup costs are fixed, but that the optimal discretionary sales decision may be different for

different demand realizations. Chan et al. [12] also incorporated the idea of tactical inventory decisions for a single-class stochastic inventory model with multi-period pricing and production decisions under limited capacity when demand is a general stochastic function. They show that when the fixed production cost is zero, then the optimal discretionary sales (or reserve inventory) is independent of the demand realization. However, when pricing is a decision, then the discretionary sales decision does depend on the demand realization.

The use of tactical inventory was extended in Chapter 2 to allow the reserving of inventory as well as planned backlogging of current customers as a second kind of tactical inventory decision. In that chapter, we consider multiple customer classes differentiated by their priority level, where the first-class customers receive complete priority over the second-class customers in the use of current resources and future backlogging. A key feature of Chapter 2 is that customers in both of the classes behave homogeneously in terms of the delivery time (all customers are willing to wait for fulfillment). The main result is that policies of the  $(S, R, B)$  form are optimal, where  $S$  is the order-up-to quantity,  $R$  is the reserve-up-to amount to protect from selling to current customers, and  $B$  is the backlog-up-to amount. Since first and second-class customers can receive delayed fulfillment, the  $R$  and  $B$  decisions may further be nested by customer class.

A fundamental difference in the current chapter compared to Chapter 2 is that in this chapter the customer classes are differentiated according to their tolerance for delay fulfillment, or “patience”. In the current chapter, customer classes are not ordered by priority on resources. Although the proof techniques in this chapter are similar, the results do not immediately follow from the models and analysis in Chapter 2 because the use of patient and impatient customers changes the form of the models. Customers differentiated according to their service preferences may be more applicable in certain settings, such as when some customer types may have an immediate need for some products.

Other papers that consider serving multiple customer classes in a production system include Deshpande et al. [21], Frank et al. [32], Gupta and Wang [38], and Sobel and Zhang [73]. In most of these, the customer classes differ according to their priority or fulfillment, and the authors look for policies to manage the system. However, a significant

difference from ours is that tactical inventory decisions are not considered in these papers. Although allowing tactical inventory may complicate the decision, previous work has shown that it can add to the profits in a manufacturing environment by providing the ability to shift demand ([12]).

We focus on a single product sold at a single manufacturer over a multi-period time horizon, where the manufacturer has limited production capacity. First-class customers claim the item immediately and never accept a delayed fulfillment and are willing to pay a premium over the market price. Second-class customers are sensitive to price, always pay the market price, and accept a delay fulfillment. We analyze the system where the manufacturer can differentiate between the customer classes and service-discriminates according to their preferences, and we also study a system where the manufacturer cannot discriminate and does not serve the classes differently though customers accept or not according to their service preferences. For both systems, the manufacturer decides in each period the amount of inventory to protect from being sold to the current period's demand and saved for future demand, and the amount of demand to backlog as well as the overall production quantity. We show that a modified base-stock policy in the form of  $(S, R, B)$  is optimal, whether the manufacturer can or cannot differentiate between the customer classes.

## ***3.2 Models and Results***

### **3.2.1 Assumptions and Notation**

We study a multi-period time horizon with periodic review where the periods are denoted as  $t = 1, 2, \dots, T$ , with  $T$  being the end of the horizon. The production in each period  $t$  is limited by the capacity,  $q_t$ , and the manufacturer pays a production cost per unit of  $c_t$  in period  $t$ . The salvage value of any units left at the end of the horizon is  $v$ . The inventory holding cost per unit in period  $t$ ,  $h_t$ , is assessed to carry inventory from period  $t$  to  $t + 1$ .

The first-class customers (index of 1) are willing to pay a premium over the market price for immediate delivery of the item and do not accept delayed fulfillment. The second-class customers (index of 2) pay the market price, and they accept fulfillment delayed up to one period; they can also be served before that deadline if resources are available. We assume

that the firm has predetermined prices  $p_t^1$  and  $p_t^2$  to charge customer class 1 and 2 in period  $t$ , respectively; the prices may be varying from period to period and are unknown to the customers until the beginning of the period  $t$ , since in some companies pricing decisions are made by the marketing department before the start of a selling season while production decisions are made by the operations department. In period  $t$ , for  $i = 1, 2$ ,  $\ell_t^i$  is the penalty per unit for demand in class  $i$  that is not satisfied and lost, and  $\beta_t^2$  is the penalty per unit for items that are backlogged for delayed fulfillment. Penalty terms for the first class are assumed to be higher than the second class.

We assume the demand of each class  $i$  in time period  $t$ ,  $D_t^i$ , is a non-stationary stochastic function; the probability and cumulative distribution functions  $(\phi_t^i, \Phi_t^i)$  are known, continuous and differentiable; and that the customer classes are independent. We do not assume particular forms of the demand functions.

The net inventory (*on-hand inventory – backlogs*) at the beginning of period  $t$  is  $I_t$ . At the beginning of a period, the manufacturer checks the inventory level and decides the production quantity; let  $S_t$  represent the inventory plus production in period  $t$ . We assume products arrive immediately, and the manufacturer fulfills the backorders carried from the previous period with the available inventory. (We allow backordered items to be delivered no more than one period later, and we restrict backorders in each period to be no more than the capacity in the next period; therefore, previously accepted orders are fulfilled before new orders are accepted.) Then the demand is realized during the period, and at the end of the period the manufacturer decides the amount of inventory to reserve for future sales and the amount of backorders to be promised in the current period for future fulfillment. Then the current demand is satisfied according to the inventory and backlogging decisions.

Let  $J_t(I_t)$  be the expected profit from period  $t$  forward to the end of the horizon when starting at period  $t$  with  $I_t$  units in inventory, or the *profit-to-go* function. Let  $G_t(S_t)$  be the expected profit-to-go from period  $t$  forward to the end of the horizon with  $S_t$  units of product available (after production). The first and second derivatives of  $J_t(I_t)$  are denoted by  $J_t'(I_t)$  and  $J_t''(I_t)$ , respectively. When the expected profit functions are specifically defined for a strategy, they will have an additional superscript indicating the strategy.

**Table 4:** Additional notation

$S_t^2$	$= (S_t - R_t^1 - D_t^1)^+$	available inventory after first-class demand is satisfied
$A_t$	$= \min(B_t, [D_t^2 - [S_t^2 - R_t^2]^+]^+)$	actual backlogged orders
$I_{t+1}^{R^1}$	$= \min(S_t, R_t^1)$	inventory carried forward due to $R^1$ decision
$I_{t+1}^{R^2}$	$= \min(S_t^2, R_t^2)$	inventory carried forward due to $R^2$ decision
$I_{t+1}^{low}$	$= [S_t^2 - R_t^2 - D_t^2]^+$	inventory carried forward due to low demand

### 3.2.2 Time Differentiation Strategy

In the Time Differentiation Strategy (TDS), we assume that the manufacturer can differentiate the customer classes by offering two time-differentiated services: selling the item for  $p_t^1$  and delivering the item immediately, or selling the item for the discounted price  $p_t^2$  and delivering the item no later than one period later.

Let  $B_t$  be the maximum planned backorders in period  $t$  to fulfill from future capacity. In period  $t$ ,  $R_t^1$  is the inventory to protect from being sold to first and second-class customers, and  $R_t^2$  is the additional inventory to protect from being sold to second-class customers; thus, the total amount to protect from class 2 in period  $t$  is  $R_t^1 + R_t^2$ . For convenience, define  $S_t^2$  to be the amount of inventory available to the second-class customers in period  $t$ , the actual backlogged orders from second-class customers after all demand is satisfied in period  $t$  from available inventory as  $A_t$ , and inventory carried forward to period  $t + 1$  due to reserved amounts  $R^1$  and  $R^2$  as  $I_{t+1}^{R^1}$  and  $I_{t+1}^{R^2}$ , respectively. If total demand is sufficiently *low* so that there is leftover inventory at the end of period  $t$ , we denote this inventory as  $I_{t+1}^{low}$ . See Table 4 for a summary of the additional notation.

In each period, after the demand for both classes of customers are revealed, the maximum number of first-class customers is satisfied from the available inventory on hand immediately<sup>1</sup>, and the second-class customers are satisfied from the available inventory left over after the first-class demand is satisfied and from the available backlog amounts<sup>2</sup>.

We model this resource allocation problem with service-differentiated customers as a

<sup>1</sup>The usage of on-hand inventory for a second-class demand instead of a first-class one is obviously sub-optimal.

<sup>2</sup>The second-class customers may be satisfied immediately if there is inventory, since it avoids the inventory holding cost and backloging penalty for those customers.

Markov decision process, where the state of the system is represented by the net inventory. We can now write the optimal expected profit in period  $t$  for the Time Differentiation Strategy as the following recursive equation:

$$J_t^{TDS}(I_t) = \max_{S_t: \max(0, I_t) \leq S_t \leq I_t + q_t} \left\{ -c_t(S_t - I_t) + G_t^{TDS}(S_t) \right\}, \quad \text{and} \quad (3)$$

$$G_t^{TDS}(S_t) = \max_{B_t: B_t \leq q_{t+1}; R_t^1, R_t^2: R_t^1 + R_t^2 \leq S_t} g_t^{TDS}(S_t, R_t^1, R_t^2, B_t), \quad \text{where} \quad (4)$$

$$\begin{aligned} g_t^{TDS}(S_t, R_t^1, R_t^2, B_t) = & \iint \left\{ p_t^1 \min(D_t^1, S_t - R_t^1) + p_t^2 \min(D_t^2, [S_t^2 - R_t^2]^+ + B_t) \right. \\ & - h_t(S_t^2 - R_t^2 - D_t^2)^+ - h_t \min(S_t^2, R_t^2) - h_t R_t^1 \\ & - \ell_t^1(D_t^1 - S_t + R_t^1)^+ - \ell_t^2(D_t^2 - [S_t^2 - R_t^2]^+ - B_t)^+ \\ & - \beta_t^2 \min([D_t^2 - [S_t^2 - R_t^2]^+]^+, B_t) \\ & \left. + J_{t+1}^{TDS}((I_{t+1}^{low} + I_{t+1}^{R^1} + I_{t+1}^{R^2}) - A_t) \right\} d\Phi_t^1(D_t^1) d\Phi_t^2(D_t^2). \quad (5) \end{aligned}$$

Equation (3) includes production cost and the remaining profit-to-go after production (Equation (4)), which is maximized over the tactical inventory ( $R_t^1, R_t^2$ ) and backlogging ( $B_t$ ) decisions. The first terms in Equation (5) include the revenue from first-class customers from the available inventory and the revenue from the second-class demand from available inventory and planned backlogging. The second line (third, fourth and fifth terms) includes the holding cost for leftover inventory and for the two reserving inventory decisions. The sixth and seventh terms, respectively, are the rejection penalties for unsatisfied first and second-class demand, and the eighth term is the delay penalty for backlogged demand. The last term in the equation is the profit-to-go in future periods, as a function of any leftover physical inventory and backlogged orders. For the last period of the horizon ( $T$ ), the final term is replaced by  $v(S_T - D_T^1 - D_T^2)^+$ , which includes the salvage cost for the leftover inventory. The constraints ensure that the manufacturer does not backlog more future capacity than he has in the next period or reserve more inventory than is available.

To simplify the TDS problem, we show that in an optimal policy, in every period either the amount of inventory protected from the second-class customers or the amount of backlogged demand of the second-class demand must equal zero (or both).

**Lemma 3.1.** *In any optimal policy under the Time Differentiation Strategy, we have:*

$$B_t \cdot (R_t^1 + R_t^2) = 0 \quad t = 1, 2, \dots, T.$$

See the Appendix for details about proofs. To see the result intuitively, suppose that  $(R_t^1 + R_t^2) > 0$ , which means that the manufacturer may reject some current second-class demand in period  $t$  in order to reserve some inventory for period  $t + 1$ . Then it is intuitive that it would not be optimal for the manufacturer to use the inventory in period  $t + 1$  to fulfill any current second-class demand in period  $t$ , thus we will have  $B_t = 0$ . The intuitive explanation for the case with  $B_t > 0$  is similar.

Lemma 3.1 implies that the structure of the optimal policies can be simplified as follows. In each period, the manufacturer can choose one of two policies: either the Reserve-Inventory policy with  $R_t^1 + R_t^2 \geq 0$ , or the Backlog-Demand policy with  $B_t \geq 0$ . Thus,

$$G_t^{TDS}(S_t) = \max\left\{G_t^{TDS-R}(S_t), G_t^{TDS-B}(S_t)\right\},$$

where  $G_t^{TDS-R}(S_t)$ , and  $G_t^{TDS-B}(S_t)$  represent the profit-to-go with  $S_t$  units of products available after production under the Reserve-Inventory policy and the Backlog-Demand policy, respectively. These policies are defined by

$$G_t^{TDS-R}(S_t) = \max_{R_t^1, R_t^2: R_t^1 + R_t^2 \leq S_t} \left\{g_t^{TDS}(S_t, R_t^1, R_t^2, 0)\right\} \quad \text{and}$$

$$G_t^{TDS-B}(S_t) = \max_{B_t: B_t \leq q_{t+1}} \left\{g_t^{TDS}(S_t, 0, 0, B_t)\right\}.$$

We will address the structural results and corresponding policies for these models in Section 3.3, after introducing the non-differentiating strategy.

### 3.2.3 Common Service Strategy

In some cases, even though the manufacturer knows the existence of multiple classes of customers, he may not be able or willing to treat customers differently. In such environments, the manufacturer manages the customers as a single class, and attempts to serve each customer with the same service strategy. We model the problem of a manufacturer who does not differentiate between two classes of customers with the *Common Service Strategy* (CSS), where the manufacturer serves customers with a first-come-first-serve rule and offers

all customers a one-period backlog for the item if the on-hand inventory is depleted. The second-class customers will accept the delayed fulfillment, but the first-class demand is lost if it is not fulfilled immediately; all customers are willing to accept immediate fulfillment.

We assume the manufacturer takes the second-class customers' reservation price,  $p_t^2$ , as the selling price to all customers, although our results also hold under other prices. The total amount of demand (class 1 and class 2 together) in period  $t$  is  $D_t^{1,2} = D_t^1 + D_t^2$  and the total demand has the probability and cumulative distribution functions  $(\phi_t^{1,2}, \Phi_t^{1,2})$ . We let  $\alpha_t$  be the average proportion of demand from the second class in period  $t$ , i.e.,  $\alpha_t = E[D_t^2]/E[D_t^1 + D_t^2]$ . We assume that the customer classes are distributed homogeneously across a time period in accordance with  $\alpha_t$ .<sup>3</sup> Let  $\ell_t$  be rejection penalty in a period  $t$ ; e.g., in our calculations we use  $\ell_t$  as the weighted average rejection penalty ( $\ell_t = (1 - \alpha_t)\ell_t^1 + \alpha_t\ell_t^2$ ), but other values can also be used. In each period  $t$ , the manufacturer decides  $B_t$ , the amount of planned backlogging in the current period;  $R_t$ , the inventory to protect from being sold in the current period; and  $S_t$ , the target level of inventory. The optimal decisions are found by solving the profit-to-go function under the Common Service Strategy:

$$J_t^{CSS}(I_t) = \max_{S_t: \max(0, I_t) \leq S_t \leq I_t + q_t} \left\{ -c_t(S_t - I_t) + G_t^{CSS}(S_t) \right\}, \text{ and} \quad (6)$$

$$G_t^{CSS}(S_t) = \max_{R_t: R_t \leq S_t; B_t: B_t \leq q_{t+1}} g_t^{CSS}(S_t, R_t, B_t) \quad \text{where,} \quad (7)$$

$$\begin{aligned} g_t^{CSS}(S_t, R_t, B_t) = & \int \left\{ p_t^2 \min(D_t^{1,2}, S_t - R_t + \min(B_t, \alpha_t(D_t^{1,2} - S_t + R_t)^+)) \right. \\ & - h_t \max(R_t, S_t - D_t^{1,2}) - \beta_t^2 \min(B_t, \alpha_t(D_t^{1,2} - S_t + R_t)^+) \\ & - \ell_t^1(1 - \alpha_t) \min(B_t/\alpha_t, (D_t^{1,2} - S_t + R_t)^+) - \ell_t(D_t^{1,2} - S_t + R_t - B_t/\alpha_t)^+ \\ & \left. + J_{t+1}^{CSS}(\max(R_t, S_t - D_t^{1,2}) - \min(\alpha_t(D_t^{1,2} - S_t + R_t)^+, B_t)) \right\} d\Phi_t^{1,2}(D_t^{1,2}) \end{aligned} \quad (8)$$

Equations (6) and (7) are as described before, except in the CSS strategy the latter is optimized over fewer reserving decisions; other differences are as below. The selling revenue includes items from both classes sold immediately, and any items backlogged from the second class only. Delay penalties are charged for backlogged second-class demand, and penalties are paid for first-class demand not satisfied immediately. The fourth term is the

---

<sup>3</sup>Note that if this assumption does not hold, the model becomes an approximation of the true situation.



penalty associated with the lost first-class demand who are offered delayed fulfillment but are not willing to accept it, and the fifth term is the rejection penalty for demand beyond the acceptance level for both classes. The last term in Equation (8) is again the profit-to-go, and constraints are as before.

Next we show that in an optimal policy, in any period, either the amount of reserved inventory equals zero or the amount of backlogged demand equals zero, i.e., they cannot both be positive.

**Lemma 3.2.** *In any optimal policy under the Common Service Strategy, we have  $R_t \cdot B_t = 0$ , for  $t = 1, 2, \dots, T$ .*

This is similar to the result for TDS, except now it applies to the reserving decision that is common to the two customer classes.

As before, with Lemma 3.2, the structure of the optimal policies can be simplified. Under the Common Service Strategy, in any period the manufacturer chooses one of two policies: either he protects inventory for the future and does not backlog current demand ( $R_t \geq 0, B_t = 0$ ), called the Reserve-Inventory policy, or he backlogs current demand but does not save items for the future ( $B_t \geq 0, R_t = 0$ ), called the Backlog-Demand policy. Thus,

$$G_t^{CSS}(S_t) = \max\{G_t^{CSS-R}(S_t), G_t^{CSS-B}(S_t)\}, \quad \text{where}$$

$$G_t^{CSS-R}(S_t) = \max_{R_t: 0 \leq R_t \leq S_t} \{g_t^{CSS}(S_t, R_t, 0)\} \quad \text{and} \quad G_t^{CSS-B}(S_t) = \max_{B_t: 0 \leq B_t \leq q_{t+1}} \{g_t^{CSS}(S_t, 0, B_t)\}.$$

There are similarities in the structure of the results for TDS and CSS, although the models have several important differences. In the next section we further analyze similarities in the structure.

### 3.3 Results

Under both the Time Differentiation Strategy and the Common Service Strategy, we can show that the four profit-to-go functions  $g_t^{TDS}(S_t, R_t^1, R_t^2, 0)$ ,  $g_t^{TDS}(S_t, 0, 0, B_t)$ ,  $g_t^{CSS}(S_t, R_t, 0)$ , and  $g_t^{CSS}(S_t, 0, B_t)$  are quasi-concave, each of them has a unique unconstrained optimizer that is independent of the inventory level  $S_t$ , and the expected profit  $J_t(I_t)$  and  $G_t(S_t)$  are

concave functions of inventory  $I_t$  and  $S_t$  respectively. These results are summarized in the following theorem:

**Theorem 3.1.** *For all  $t = 1, \dots, T$ ,*

- $g_t^{TDS}(S_t, R_t^1, R_t^2, 0)$  is a jointly quasi-concave function of  $R_t^1$  and  $R_t^2$ , and  $g_t^{CSS}(S_t, R_t, 0)$  is a quasi-concave function of  $R_t$ ,
- $g_t^{TDS}(S_t, 0, 0, B_t)$  and  $g_t^{CSS}(S_t, 0, B_t)$  are quasi-concave functions of  $B_t$ ,
- $G_t^{TDS}(S_t)$  and  $G_t^{CSS}(S_t)$  are concave functions of  $S_t$ ,
- $J_t^{TDS}(I_t)$  and  $J_t^{CSS}(I_t)$  are concave functions of  $I_t$ ,
- The unconstrained optimizers  $(R_t^{1*}, R_t^{2*}, \text{ and } B_t^*)$  for functions  $g_t^{TDS}(S_t, R_t^1, R_t^2, 0)$ , and  $g_t^{TDS}(S_t, 0, 0, B_t)$ , are independent of inventory level  $S_t$ , where for  $i = 1, 2$ ,

$$(R_t^{1*}(S_t), R_t^{2*}(S_t)) = \operatorname{argmax}_{(R_t^1, R_t^2): 0 \leq R_t^1, 0 \leq R_t^2} \left\{ g_t^{TDS}(S_t, R_t^1, R_t^2, 0) \right\} \text{ and}$$

$$B_t^*(S_t) = \operatorname{argmax}_{B_t: 0 \leq B_t} \left\{ g_t^{TDS}(S_t, 0, 0, B_t) \right\}.$$

- The unconstrained optimizers  $(R_t^*$  and  $B_t^*)$  for  $g_t^{CSS}(S_t, R_t, 0)$  and  $g_t^{CSS}(S_t, 0, B_t)$  are independent of inventory level  $S_t$ , where

$$R_t^*(S_t) = \operatorname{argmax}_{R_t: 0 \leq R_t} \left\{ g_t^{CSS}(S_t, R_t, 0) \right\}, B_t^*(S_t) = \operatorname{argmax}_{B_t: 0 \leq B_t} \left\{ g_t^{CSS}(S_t, 0, B_t) \right\}.$$

Theorem 3.1 implies an optimal policy for both the Time Differentiation Strategy and the Common Service Strategy that has a similar form, and thus we have the following corollary.

**Corollary 3.1.** *Given a vector of prices, there exists an optimal policy for*

- the Time Differentiation Strategy with an optimal order-up-to level  $(S_t^*)$ , optimal reserve-up-to-levels  $(R_t^{1*}$  and  $R_t^{2*})$ , and an optimal backlog-up-to level  $(B_t^*)$ ,
- the Common Service Strategy with an optimal order-up-to level  $(S_t^*)$ , an optimal reserve-up-to-level  $(R_t^*)$  and an optimal backlog-up-to level  $(B_t^*)$ .

Note that for CSS there is a single reserve inventory decision that non-discriminatingly applies to both classes of customers, and similarly for the planned backlogging decisions, while for TDS there are separate values for reserving that apply to each class and the planned backlogging only applies to the second-class demand. However, in both cases the form of the optimal policy is  $(S, R, B)$ . In both cases the optimal policies are considered to be modified base stock ones, because the realized values may be limited by capacity or available inventory. The results also show that the optimal inventory decisions are independent of the realized demand, which implies that the decisions could also have been made before the exact demand realization.

Examining the decisions in more detail provides more information on their meaning. The optimal decisions for CSS are defined by the following:

$$\begin{aligned}
S_t^* &= \max\{S : c_t \leq G_t'^{CSS}(S)\} && \text{if } c_t \leq G_t'^{CSS}(0) \\
R_t^* &= \max\{I : p_t^2 + \ell_t + h_t \leq J_{t+1}'^{CSS}(I)\} && \text{if } p_t^2 + \ell_t + h_t < J_{t+1}'^{CSS}(0) \\
B_t^* &= \min\{I : J_{t+1}'^{CSS}(-I) \geq p_t^2 + \ell_t^2 - \beta_t^2\} && \text{if } p_t^2 + \ell_t^2 - \beta_t^2 > J_{t+1}'^{CSS}(0),
\end{aligned} \tag{9}$$

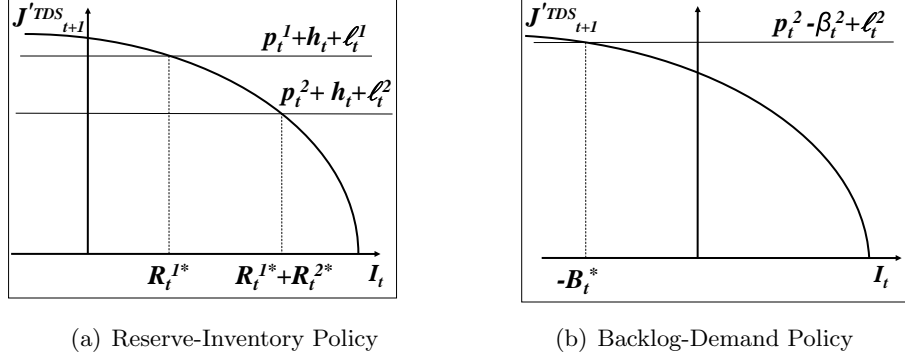
and the optimal decisions for TDS that are different from CSS are given by:

$$\begin{aligned}
R_t^{1*} &= \max\{I : p_t^1 + \ell_t^1 + h_t \leq J_{t+1}'^{TDS}(I)\} && \text{if } p_t^1 + \ell_t^1 + h_t < J_{t+1}'^{TDS}(0) \\
R_t^{1*} + R_t^{2*} &= \max\{I : p_t^2 + \ell_t^2 + h_t \leq J_{t+1}'^{TDS}(I)\} && \text{if } p_t^2 + \ell_t^2 + h_t < J_{t+1}'^{TDS}(0).
\end{aligned} \tag{10}$$

For each decision, if the condition is not satisfied, then the decision variable equals zero.

In each case, a decision is found by comparing the net revenue from gaining or losing a customer with the marginal expected profit of an additional unit in the future ( $p_t^i + \ell_t^i + h_t$  is the net revenue of selling to customer class  $i$  from inventory,  $p_t^i + \ell_t^i - \beta_t^i$  is the net revenue from backlogging an  $i$  class customer, and  $J'(I)$  is the marginal expected profit of an additional unit above  $I$ ). This is also apparent from examining the decisions for TDS that are pictured in Figure 5 (the ones for CSS are similar in structure).

Chapter 2 examined the form of the optimal policies when customers are differentiated by priority level and behave homogeneously with respect to delayed fulfillment. In that case, the form of the decisions was similar (i.e.,  $(S, R, B)$ ), but a different set of decisions



**Figure 5:** Optimal Decisions for TDS under the Optimal Policies

and policies resulted from the model. An obvious difference is that in the current research only one backlogging decision results, since only the second-class customers are willing to accepted delayed fulfillment. Another difference is that in Chapter 2, one submodel that we found could be optimal was the Reserve-and-Backlog Policy, where both a backlogging and reserving decision could be positive in the current period, but that is not true for the TDS and CSS models.

Numerical experiments of the TDS and CSS policies in the next section show that significant profit improvement can be achieved with the tactical inventory, especially when production capacity is limited.

### 3.4 Computational Analysis

To further analyze the impacts of time differentiation and the corresponding tactical inventory decisions, we perform a computational study. The benchmark against which we compare TDS and CSS is a traditional production and inventory problem with limited capacity and no tactical inventory decisions with all customers served as a single class (unsatisfied demand is lost).

We compare the performance of the  $(S, R, B)$  type policies over the traditional policies using the metric of profit potential, as defined by  $100 * (\frac{V_{(S,R,B)}}{V_S} - 1)$ , where  $V$  indicates the expected profit of the problem being solved. In CSS and the traditional policy, we charge price  $p_t^2$  to all customers. This implies that TDS may show a big improvement in profit that is due, in part, to the ability to differentiate customers.

The profit improvement of the Common Service Strategy compared to the traditional inventory policy comes from a single source: the tactical inventory (reserving and planned backlogging), whereas the Time Differentiation Strategy has both tactical inventory as well as differentiated pricing. Thus, by comparing both TDS and CSS to the traditional policy, we can separate the impacts of price differentiation versus tactical inventory.

### 3.4.1 Experiment Details

In each experiment, the average total demand from first and second-class customers is 100, which we refer to as  $Dem^*$ . We assume that demand uncertainty is additive with a mean of 0. Since demand variation is usually proportional to average demand, we set the ratio between the standard deviation and the expectation to be 20% in each period.

Production capacity is constant throughout the planning horizon of an instance. Across experiments we use three values of production capacity: 60 (low), 80 (med), and 100 (high). The production cost may vary by period, but the production cost vector is the same across instances. See Table 5 for the exact data; for example, the average markup of  $p_t^2$  ( $p_t^1$ ) over the cost is about 30% (60%) for the experiments on class proportions.

We study the impact of proportion of the second class customers,  $\alpha_t$ , in our first set of experiments. We let  $\alpha_t$  take the values of 0.2, 0.25, 0.5, 0.75, 0.8 which corresponds to the expected demand from second-class customers over the horizon taking the values of 20, 25, 50, 75, and 80. The prices vary by period but are the same across all experiments. This creates the likelihood that different submodels will be optimal in a particular instance. The average ratio of  $p_t^1/p_t^2$  is 1.22 for the experiments studying the proportion of demand. See Table 5 for the prices used in this set of experiments; the data was chosen so that comparisons can be made between the models in this chapter and the models in Chapter 2.

### 3.4.2 Results

Observe that the relative performance of CSS and TDS for a given proportion of second-class customers ( $\alpha_t$ ) is better when capacity is tight. As an example, in Figure 6(a) the performance of TDS when capacity is 0.6  $Dem^*$  is better than the performance of TDS

**Table 5:** Specific experimental data

t	1	2	3	4	5	6	7	8	9	10	11	12	Avg
$c_t$	70	90	70	50	70	90	70	50	70	90	70	50	<b>70</b>
$p_t^2$	90	110	90	70	90	110	90	70	90	110	90	70	<b>90</b>
$p_t^1$	110	130	110	90	110	130	110	90	110	130	110	90	<b>110</b>

when capacity is  $0.8 Dem^*$ . The figure also shows that impact on TDS and CSS as  $\alpha_t$  increases. Not surprisingly, the relative performance reduction in TDS as  $\alpha_t$  increases can easily be explained by the loss of additional revenue opportunities from first-class customers. CSS performs better as the proportion of second-class customers increases. This is in part because more customers are willing to accept delayed fulfillment, and it may also be explained increased flexibility in the managing of demand due to backlogging.

The figure also shows the impact of tactical inventory, seen by comparing CSS to the traditional policy, since both offer  $p_t^2$  to all customers. The profit improvement is 17% or more in the figures.

In order to see the value of being able to differentiate between the customer classes, we set the price  $p_t^1$ , that is charged to the first-class customers, to  $p_t^2$  in Figure 6(b). We observe a profit improvement as high as 6% over CSS even though the same price is charged to everyone. The main insight from Figure 6(b) is that even though the price that will be charged to the first-class customers are close to the market price, there is a significant amount of profit improvement opportunity for the manufacturer from the flexibility that is gained.

Another interesting observation from the graphs is that the profit function is not linear with respect to the second-class proportion of demand. This may even suggest that if firms are able to change the mix of customers in their market, that this has the most impact when the proportion of first and second-class customers is approximately the same.

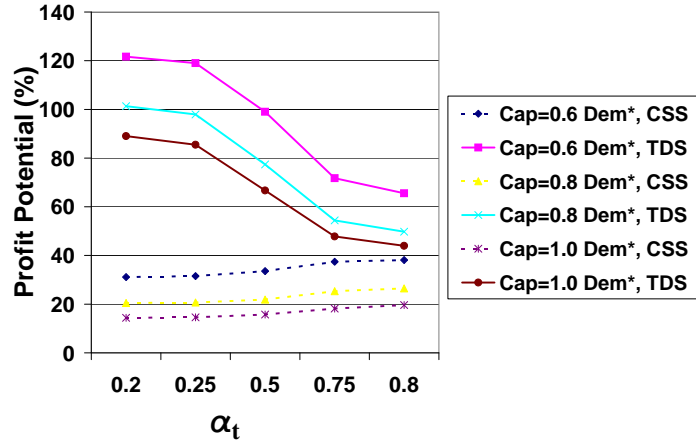
### 3.5 Conclusions

Many companies today provide differentiated service so that some customers are served immediately while others receive delayed fulfillment for a discount. Customers receive higher

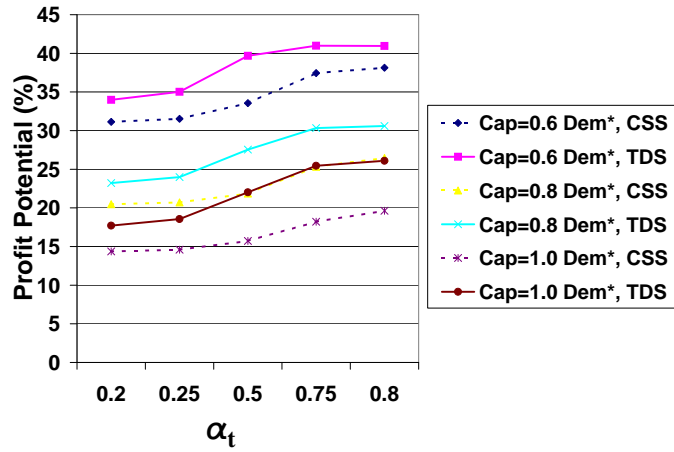
utility, and the company may gain flexibility to improve operations. We address such a situation in this chapter where the first-class customers require immediate service and pay a premium for immediate response from the firm, and the second-class customers accept delayed fulfillment and pay the market price. We investigate where the firm has the ability to differentiate the customer classes or does not (Time Differentiation Strategy and Common Service Strategy, respectively).

We consider these problems in a stochastic production and inventory context, and we use tactical and planned inventory decisions in order to allocate scarce inventory and production resources effectively. We show that for general stochastic demand functions the optimal policies for both TDS and CSS are modified base-stock policies in the form of  $(S, R, B)$ , where  $S$  is the optimal order-up-to level,  $R$  denotes the optimal reserve-up-to levels, and  $B$  is the optimal backlog-up-to levels.

This chapter contributes to the literature on operational models to manage markets with segmented demand, and it shows the impact of one kind of flexibility in the production system. Additional research would also be helpful in showing how to manage systems with segmented demand. For instance, pricing can be used to determine the size of the customer classes in response to variability, and policies to manage systems with extended leadtimes and multiple classes is another. The area is rich and has many applications in modern manufacturing and e-tailing companies that may be experimenting with different business models.



(a) Performances with price differentiation in TDS



(b) Performances when  $p_t^1 = p_t^2$  in TDS

**Figure 6:** The relative performance of  $(S, R, B)$  type policies over traditional inventory decisions



## CHAPTER IV

### LEADTIME QUOTATION AND ORDER ACCEPTANCE WHEN DEMAND DEPENDS ON SERVICE PERFORMANCE

#### *4.1 Introduction*

In this chapter we investigate a firm's leadtime quotation decisions when its recent performance of meeting quoted leadtimes affects the arrival of future orders. As customers increasingly look for better and faster service, when choosing a supplier they consider both the length and the "reliability" of the quoted leadtimes. While short leadtimes are desirable, the reliability of the quoted leadtimes (i.e., the supplier's ability to meet the promised due dates, or "service level") is equally important, especially for business customers. A late delivery from a supplier can shut down a manufacturing line, for example, costing the customer of that supplier millions of dollars. 78% of companies which operate in a just-in-time environment in the U.S. ranked delivery reliability as high priority, whereas only 25% ranked price as high priority (Billesbach et al. [9]).

Suppliers can pay high penalties for late deliveries, and these penalties usually increase with the delay. For example, Real World Components offers a 105% delivery guarantee, i.e., customers get the return price plus 5% if their order does not arrive on time (Carbone [11]). Besides the immediate monetary impact of such penalties, late deliveries can also damage the image of the supplier and reduce the arrivals of future orders. Repeat customers keep track of the firms' "service level" through various means and consider the recent delivery performance of a seller when deciding whether or not to place an order. For example, as Silicon Graphics Inc. (a leader in the three-dimensional graphics computers in the 1990's) began to lose its technological edge, it started to lose customers and revenues due to not meeting the leadtimes quoted to the customers. As stated in the cover story of Business Week Online ( Hof et al. [43]), longtime SGI loyal customers started to drop SGI since they stopped believing SGI salespeople's insistence that they would ship on time. Similarly, even

for one-time retail customers, sites such as pricegrabber.com provide merchant ratings that allow customers to access e-tailers' past delivery performance.

While recent research on leadtime decisions has considered the impact of the quoted leadtime on a customer's decision to place an order, much has ignored the impact of the seller's past performance in meeting the promised delivery dates. However, as the examples above indicate, in today's environment with easily accessible information about sellers' delivery performances, both the length and also the reliability of the quoted leadtimes impact customers' decisions in choosing a supplier. In this chapter, we model customers' sensitivity to the reliability as well as the duration of the quoted leadtimes, and we demonstrate how to quote leadtimes and determine how many orders to accept in infinite and finite capacity settings. We show that the impact of ignoring past service can be significant, including the quotation of "unethical leadtimes" or even going out of business, while incorporating service performance can increase revenue significantly. To the best of our knowledge, this is the first in the leadtime literature to consider the impact of the firm's past performance in meeting delivery promises on the arrival of future orders.

## ***4.2 Literature Review***

Most due date management policies proposed in the early literature assume that customers accept the quoted leadtimes (due dates) regardless of their duration (Baker [2], Baker and Bertrand [3], Enns [25], Fry et al. [33], Hopp and Roof Sturgis [44], Miyazaki [60], Spearman and Zhang [74], Weeks [86], Wein [87]). Many of these papers propose a two-step approach: assign the due dates first, and then schedule the orders using a priority dispatch policy such as first come first serve, shortest processing time, earliest due date, etc. A common approach for setting due dates is to use dispatch due date rules which follow the general form  $d_j = r_j + f_j$  where  $d_j$ ,  $r_j$ , and  $f_j$  are the due date, the release time, and the *flow allowance* for job  $j$ . The tightness of the flow allowances (and the due dates) is usually controlled by some parameters. To ensure that the assigned leadtimes are reliable to the extent possible, they either include a lateness penalty in the objective function or impose a service level constraint, such as the average fraction of tardy jobs or maximum expected

tardiness (e.g., see [44], [74], [87]).

In most businesses, the quoted leadtimes (or due dates) and price affect the customers' decisions to place an order. Equivalently, the firm has the choice of accepting or rejecting an order. For example, a customer with a firm deadline may not place the order if the quoted due date exceeds the deadline. In general, the longer the quoted leadtime, the less likely that a customer will place an order. Recent papers in the literature that study leadtime decisions capture the impact of the quoted leadtimes on demand, assuming that the probability that an arriving customer places an order decreases as the quoted leadtime increases (Chatterjee et al. [15], Dellaert [18], Duenyas [23], Duenyas and Hopp [24], Slotnick and Sobel [69]) or the customer does not place an order if the quoted due date exceeds the customer's deadline (Charnsirisakskul et al. [13], Keskinocak et al. [48]). Hence, due date quotation decisions are considered together with order acceptance decisions, taking a profit maximization rather than a cost minimization perspective. In general, the revenue from an accepted order (in class  $j$ ) is  $R$  ( $R_j$ ) and there are earliness/tardiness penalties if the order is completed before/after its quoted due date. Let  $P(l)$  denote the probability that a customer places an order given quoted lead time  $l$  and let  $l_{max}$  denote the maximum acceptable leadtime to the customer. The proposed demand models in the literature include the following:

$$(D1) : P(l) = 1 - l/l_{max}$$

$$(D2) : P(l) = \begin{cases} 1, & \text{if } l \leq l_{max} \\ 0, & \text{otherwise} \end{cases}$$

$$(D3) : P(l) = e^{-\lambda l}, \text{ where } \lambda \text{ is the arrival rate of the customers}$$

$$(D4) : P(l) \text{ is a decreasing concave function of } l$$

The paper that is most closely related to this chapter is [24], where the authors consider demand models (D2)-(D4). They first consider a system with infinite server capacity and for the special case of exponential processing times and model (D3), they find a closed form solution for the optimal leadtime. Next, they consider the capacitated case studying a single server queue  $GI/GI/1$  where processing times have a distribution in the form of increasing failure rate (IFR). They first study the problem for model (D2) where the firm's main decision is to decide which orders to accept (reject), by quoting a leadtime less (greater)

than  $l_{max}$ . They show that the optimal policy has a control-limit structure: for any  $n$ , the number of orders currently in the system, there exists a time  $t(n)$  such that a new order is accepted if the first order has been in service for more than  $t(n)$  time units. For model (D4) and an  $M/M/1$  queue, they show that the optimal leadtime to quote is increasing in  $n$ . [23] extends some of these results to multiple customer classes, with different net revenues and leadtime preferences.

We extend the infinite capacity model in [24] to incorporate service. It is interesting to note that when the impact of the service on future arrivals is ignored, the firm might find it profitable to quote “unethical” leadtimes, even if there is a service constraint. Spearman and Zhang [74] study the leadtime quotation problem with the objective of minimizing the weighted average leadtime subject to an upper bound on the average tardiness. They show that it is optimal to quote a zero leadtime if the congestion level is above a certain threshold, even though the possibility of meeting this quote is extremely low. Intuitively, when the system is congested, an arriving job will be late with high probability, unless a very long leadtime is quoted. However, long leadtimes negatively affect the objective function. Since the service level is on the number of tardy jobs, when the system is congested it is preferable to simply quote a zero leadtime, adding to the number of tardy jobs but keeping the objective function value low. Using numerical examples, the authors show that a customer is more likely to be quoted a zero leadtime when the service level is low or moderate, rather than high, creating service expectations completely opposite of what the system can deliver. In contrast, we show that when the firm considers the impact of service on future arrivals, it is never profitable to quote a zero leadtime.

Another stream of recent papers within the due date management literature considers due date and price decisions simultaneously (Boyaci and Ray [10], Charnsirisakskul et al. [14], Palaka et al. [62], Ray and Jewkes [64], So [71], So and Song [72]). There is also a small but growing literature considering leadtime (and price) decisions within a decentralized marketing-operations framework ([15], Pekgun et al. [63]). However, in contrast to our dynamic setting, most of these papers study queuing models focusing on a common due date in steady state.

In all the papers cited in the due date management literature above, current performance in service quality (meeting the quoted leadtimes) does not affect future customer arrivals. There are studies in the economics and operations management literature that consider the effect of service quality on the future market share. In the economics literature, Deneckere and Peck [19] report a positive correlation between price charged and service quality in a two-stage game where service quality is measured by per customer capacity. Liebeskind and Rumelt [53] consider a case where firms have the option of selecting to be a high or low-quality producer and show that the firm will be honest about his choice of service level only if the price charged to the customers is equal to their reservation price. In the operations management literature, Mendelson and Whang [58], Stidham [77] and van Mieghem [82] use queueing settings for single firms to optimize the system wide performance via price mechanisms where delay in queue or system is used as the quality measure. Hall and Porteus [41] utilize a simple dynamic model to investigate the behavior of the firm where firms compete by investing in capacity (and capacity implies the delivered service quality). Gans [36] focuses on the firm's choice of mean service quality by considering the long-run average profit, where a customer chooses a firm using the history of the service quality he received so far. Similar to [36], we also consider a single static decision that affects the service level to be consistent with the industry practice of stationary targets for service levels instead of allowing the firm to change service levels in response to short-term changes as in [41]. While [41] and [36] study (in a game theoretic setting) a closed system where  $n$  customers switch from one firm to another, we consider an open system in a queueing setting.

In this chapter we extend the literature on leadtime quotation decisions by considering *the impact of the firm's ability to meet the quoted leadtimes on future arrivals* as well as *the impact of the quoted leadtime on the probability of an arriving customer's order placement decision*. For a recent review of due date management policies, see Keskinocak and Tayur [49].

### 4.3 Models and Results

We study leadtime quotation decisions in a stochastic environment where a leadtime,  $\ell$ , is quoted to each arriving customer and the customer's probability of placing an order decreases in the quoted leadtime. We make the following modeling assumptions:

- the expected arrival rate of the customers is  $\lambda(s)$ , where  $\lambda$  can be interpreted as the base arrival rate and  $\lambda(s)$  is an increasing function of the service level;
- the service time of an order is exponentially distributed with mean  $1/\mu$ ;
- placed orders create an immediate revenue of  $R$ ;
- on-time completion of an order improves, whereas late completion of an order degrades the service level,  $s$ , of the firm; and
- if an order is not completed on-time the firm incurs a penalty,  $c$ , per unit time, i.e., the total penalty paid by the firm for a late order is proportional to the length of the delay.

We consider two decision models, namely, *Naive* and *Service-Sensitive*. In the Naive model, the firm makes leadtime decisions assuming that the arrival distribution of the customers is stationary throughout time (i.e., the arrival rate is constant and does not depend on the firm's performance in meeting the quoted leadtimes). By contrast, in the Service-Sensitive model the firm takes into account the fact that the future arrival rate of the customers is directly affected by whether or not past orders were completed on time. We study these models under infinite and finite capacity settings, analyze optimal leadtime quotation policies, and investigate the impact of the customers' service sensitivity on the firm's leadtime decisions and profits.

#### 4.3.1 Infinite Capacity Case

When the firm's capacity is unlimited (equivalently, if the firm has a high number of servers), upon placing an order each customer is assigned a server and the service of the order starts immediately. If a leadtime  $\ell$  is quoted to an arriving customer, the customer's probability of placing an order is  $p(\ell) = e^{-\alpha\ell}$ , where  $1/\alpha$  can be thought of as the mean leadtime acceptable to customers (as in [69]). We assumed that the service level  $s$  takes continuous

values between zero and one (zero indicates the lowest and one indicates the highest service level), and the expected arrival rates of the customers,  $\lambda(s)$ , is linear in service level with highest arrival rate of  $\lambda$ . The expected arrival rate of the customer is assumed to be linear in  $s$ , and its In the long-run, the service level  $s$  is the probability of meeting the quoted leadtime,  $\ell$ , where  $s = 1 - e^{-\mu\ell}$ .

If the firm is *Naive*, it ignores the impact of the service level on the arrival rate and assumes that the customers arrive at the base (expected) arrival rate  $\lambda$ . In this case the optimal leadtime is found by solving  $\Pi^N = \max_{\ell \geq 0} e^{-\alpha\ell} (R - \int_l^\infty c(y-l)\mu e^{-\mu y} dy) \lambda$  and the optimal leadtime is  $\ell^N = \frac{1}{\mu} \left( \ln \left( \frac{(\mu+\alpha)c}{\alpha R \mu} \right) \right)^+$  (see [24]). Note that if  $R \geq \frac{(\mu+\alpha)c}{\alpha\mu}$ , then  $\ell^N = 0$ .

**Observation 4.1.** *If the revenue per order is sufficiently large and the firm ignores the impact of the service level on future arrivals (or if the customers are insensitive to the service level), then this leads to the quotation of “unethical” leadtimes where every order is delivered late (see [74] for a discussion on unethical leadtimes).*

Clearly, quoting a zero leadtime to every arriving customer while there is no hope of meeting such a promise is not an ethical business practice, even if it means attracting many customers and maximizing profits in the short-term. Furthermore, what is the impact of such a behavior (or ignorance) on the firm’s long-term profitability?

To answer this question, we next consider the leadtime decisions of a *Service-Sensitive* firm. We find the optimal leadtime to be quoted by the firm by solving the following profit maximization problem:

$$\Pi = \max_{\ell \geq 0} p(\ell) \left( R - \int_\ell^\infty c(y-l)\mu e^{-\mu y} dy \right) (1 - e^{-\mu\ell}) \lambda. \quad (11)$$

**Theorem 4.1.** *For a Service-Sensitive firm,*

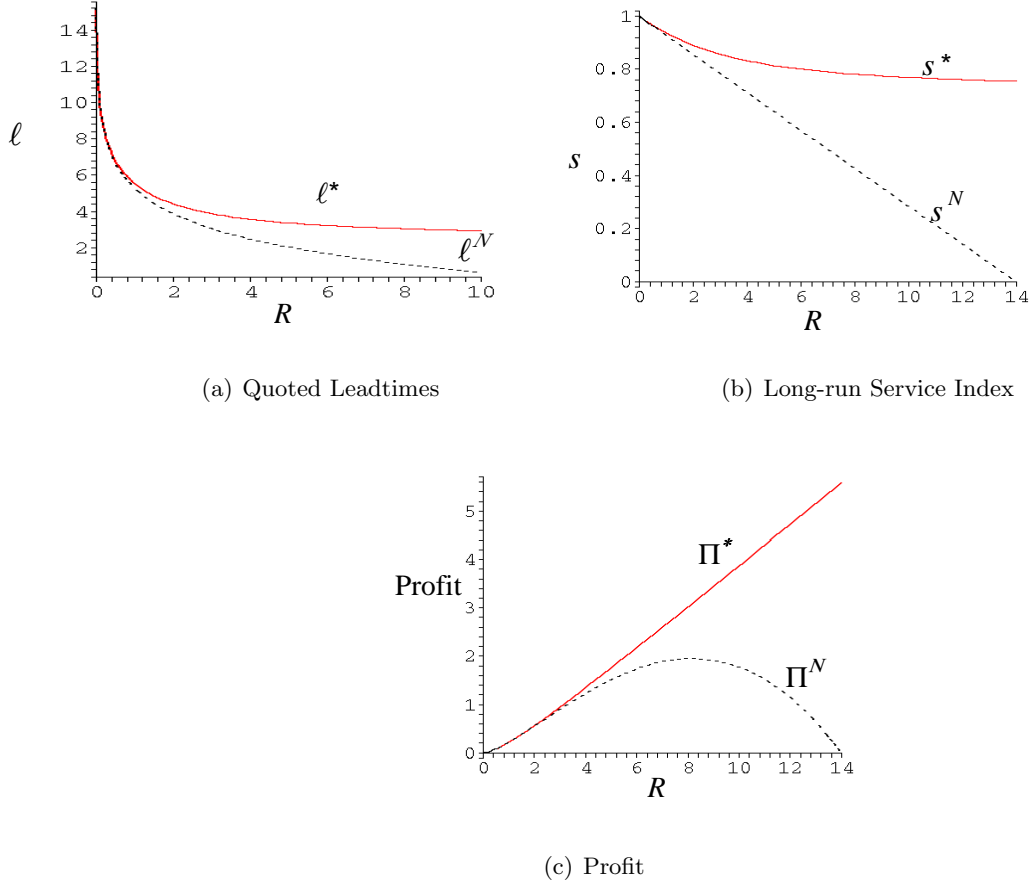
- (i) *the optimal leadtime to quote is  $\ell^* = \frac{1}{\mu} (\ln(\theta))^+$ , where  $x^+ = \max\{0, x\}$  and  $\theta = \frac{(\alpha+\mu)(c+\mu R) + \sqrt{(\alpha+\mu)^2(c-\mu R)^2 + 4\mu^3 R c}}{2\alpha R \mu}$ ;*
- (ii)  *$\ell^*$  is decreasing in  $R$  and  $\lim_{R \rightarrow \infty} \ell^* = \frac{1}{\mu} \ln\left(\frac{\alpha+\mu}{\alpha}\right) > 0$ ;*
- (iii)  *$\ell^* > \ell^N$ .*

The proofs of all the results are presented in the Appendix. From Theorem 4.1, a Service-Sensitive firm quotes strictly positive leadtimes, which are always longer (i.e., more conservative) than the ones quoted by the Naive firm. While such longer leadtimes might decrease the number of orders placed upon arrival, serving fewer orders increases the firm's service level and positively impacts the arrival stream in the long term. Hence, the Service-Sensitive firm projects a more reliable (or ethical) image to the customers in terms of meeting its delivery promises.

In Figure 7, we show a numerical example that compares the profits of the Naive and Service-Sensitive firms. From Figure 7(a) and (b), as the revenue per order ( $R$ ) increases, the deviation from the optimal leadtime increases while the service level significantly decreases (and eventually reaches zero) if the firm ignores the customers' sensitivity to the service performance. From Figure 7(c), the firm's profit increases in  $R$  when the optimal leadtime  $\ell^*$  is quoted; however, ignoring customers' sensitivity to the service performance leads to a significant loss of profit, and for large  $R$  values, causes the firm to lose all of its customers and go out of business. There are many business cases which fit this scenario. For example, consider eToys.com. According to USA Today (Krantz [50]) "Last year [1999], eToys and others were beset with delivery snafus. That hurt holiday sales this year." eToys made unrealistic delivery promises in 1999, which it could not meet, and infuriated customers who swore to never purchase from eToys again (as stated in online merchant reviews [59]). eToys closed its doors in 2001.

One might question whether it is possible to achieve the efficient system outcome (that is, incorporating customer behavior) by modifying the penalty cost ( $c$ ). Figures 8(a) and 2(b) show the effect of  $c$  on the long-run service level and profit for both a Naive and a Service-Sensitive firm. One has to charge a very high penalty cost (to include both the immediate penalty and reputation or word-of-mouth) to capture some of the overall effects of the customers' service sensitivity. For instance, when the penalty cost is 30%, 50%, 75%, 100% of the revenue, the profit and the service level of the Service-Sensitive firm are 162%, 19%, 6%, 3% and 257%, 33%, 11%, 5% higher than that of the Naive firm, respectively. Note that the higher service levels of the Service-Sensitive firm may have additional implications

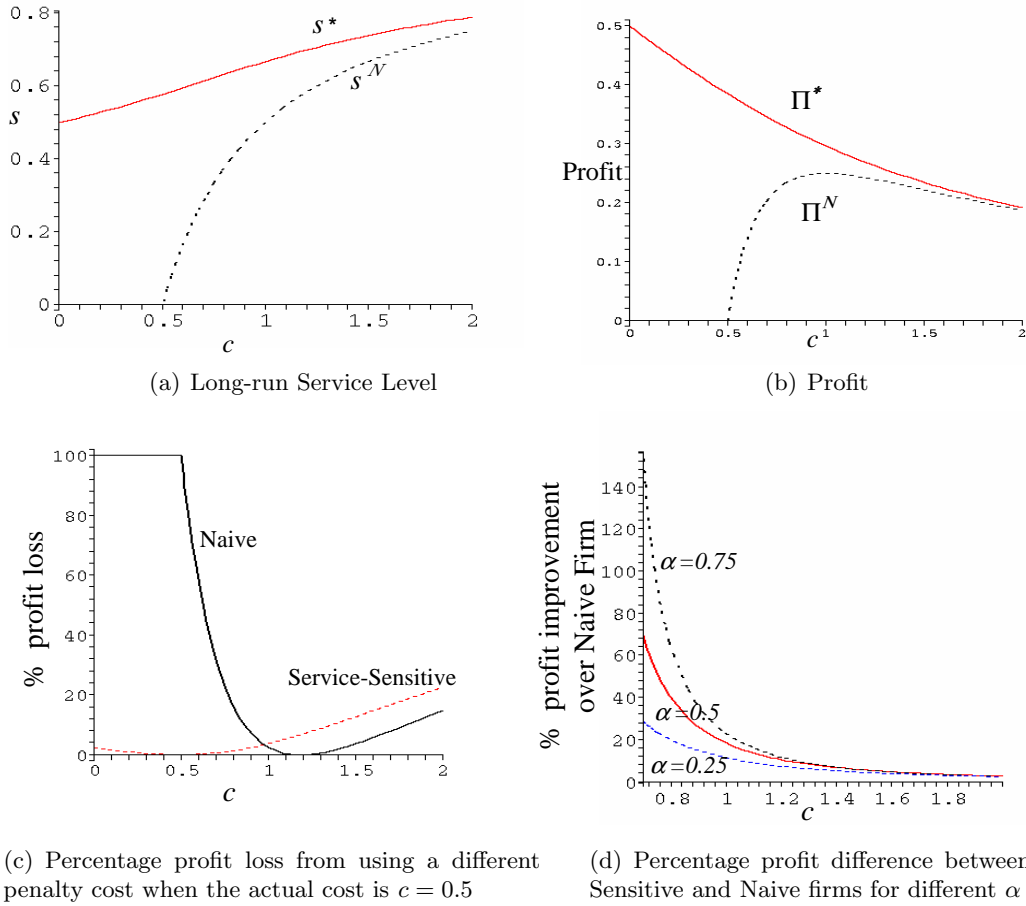




**Figure 7:** Performance based on revenue per item when  $\alpha = 0.2, \mu = 0.5, c = 2$

beyond impacting future customer arrivals; for example, they might increase the customers' willingness to pay leading to higher revenues (and profits).

It is also useful to consider how the Naive and Service-Sensitive firms perform when they make a mistake in estimating the penalty cost  $c$ . Figure 8(c) shows the percentage profit loss if a different cost than the actual  $c$  is used by either firm. Underestimating the penalty cost (or even using the actual) is detrimental to the Naive firm's profits, while the Service-Sensitive firm's profits are quite robust for a large range ( $\pm 100\%$ ) of penalty cost values around the actual. Significantly overestimating the penalty cost leads to a slightly



**Figure 8:** Performance based on penalty cost when  $\alpha = 0.5, \mu = 0.5, R = 2$

higher profit loss for the Service-Sensitive firm than the Naive firm, which is expected since the overestimation acts as a “correction” for the Naive firm’s ignorance of the customers’ service sensitivity. Overall, the profit of the Naive firm is significantly more sensitive to the choice of  $c$  than that of the Service-Sensitive firm.

A key factor in a firm’s performance can be its “relative” capacity, which depends on the mean leadtime ( $1/\alpha$ ) acceptable to customers and the mean service time ( $1/\mu$ ) (note that a higher  $\alpha$  value corresponds to the customers’ desire for shorter leadtimes, and hence, effectively lower capacity). Figure 8(d) illustrates that the percentage profit difference between Service-Sensitive and Naive firms is higher when the capacity is tight ( $1/\alpha < 1/\mu$ ) and significant even when  $c$  is close to  $R$ . For example, for  $c = 0.5R$ , the percentage profit difference is 12% and 23% under “high” and “low” capacity, respectively.

### 4.3.2 Finite Capacity Case

To study the leadtime quotation decisions under limited capacity, we consider a single server, where the service discipline is first-come-first-serve (FCFS). We assume that there is an industry standard leadtime,  $\ell$ , which is acceptable to customers (as described in [24]). That is, if the firm quotes a leadtime  $\ell$  (or less), the customer places an order, and if the quoted leadtime is larger than  $\ell$ , the customer leaves the system. Since a customer's order placement decision is the same for all leadtimes less than  $\ell$ , to avoid lateness penalties the firm will never quote a leadtime strictly less than  $\ell$ . Hence, the firm's decision is whether to accept a customer (by quoting  $\ell$ ) or to reject (by quoting longer than  $\ell$ ). Industries might have such fixed leadtimes if there is significant competition that has driven leadtime to a common value, if leadtime is small compared to the known transportation time, or if there is a batch process production. Retail firms may also have such fixed leadtimes, e.g., one hour photo processing or thirty minute pizza delivery.

As in the case of infinite capacity, we study the leadtime quotation decisions of *Naive* and *Service-Sensitive* firms. In the case of a Service-Sensitive firm, to compute the backlog when a customer arrives to the system, one needs to keep track of the remaining time until the due date for each customer (contrary to [69]), which makes the typical semi-Markov decision process model intractable. Therefore, we consider a simplified discrete model, where we set the length of a period to the industry dictated leadtime,  $\ell$ , and customers arrive in batches at the beginning of each period. Let  $\{F_s(k) : s \in S\}$  be a collection of stochastically increasing distribution functions where  $S = \{s_{min}, \dots, s_{max}\}$  denotes the set of service levels. At the beginning of a period with service level  $s$ ,  $dF_s(k)$  is the probability of having  $k$  customer arrivals. We assume  $k \leq \bar{k}$ , i.e., the number of arrivals in any period is finite. The firm accepts  $a$  customers from this batch and the orders of the accepted customers are due at the end of that period. Hence, the length of each decision epoch is also  $\ell$ . In each period, the service of  $X(\ell)$  customers is completed. The firm pays a lateness penalty of  $c$  for each customer where service is not completed at the end of the period, and these customers are carried to the next period. Customer dissatisfaction can be generally associated with worst-case performance as stated in Fleisch and Powell [31], therefore as a

proxy for estimating the service performance in this complex system, we let the service level change depending on whether the order of the last customer in the system was completed on time. Since orders are completed during discrete intervals, the service level ( $s$ ) goes up ( $s^+$ ) if all the orders in a given period are completed before the end of the period, and goes down ( $s^-$ ), otherwise. This model can be motivated by industries where orders can be placed at any time, but the state of the system is assessed periodically (e.g., on a rolling horizon basis) to update decisions. Discrete periods are desirable in any system where monitoring costs are too high for continuous review, e.g., where production plans may be determined daily or weekly.

#### 4.3.2.1 Leadtime Decisions of the Naive Firm

The *Naive* firm believes that the arrival distribution of the customers is stationary throughout time, and is independent of the performance in meeting the quoted leadtimes. Therefore, it assumes that the arrival distribution  $F(k)$  is independent of the service level  $s$  and  $dF(k)$  is the probability of having  $k$  arrivals in a given period.

We model this problem by a Markov decision process (MDP). The state of the system is represented by  $i$ , where  $i$  is the number of orders in the system at the beginning of a time period before the customers arrive. Let  $V_n(i)$  be the expected net benefit of the system over periods  $n$  ( $n \geq 0$ ) to 0 and  $V_0(i)$  be the expected reward to complete the service of the  $i$  customers remaining in the last period of the horizon. Therefore,  $V_0(i) = -[c \cdot i + E_{clearing}(i)]$ , for  $i = 0, 1, \dots$ , where,

$$E_{clearing}(i) = \frac{Pr\{X(\ell) = 0\} \cdot c \cdot i + \sum_{j=1}^i Pr\{X(\ell) = j\} [(i-j) \cdot c + E_{clearing}(i-j)]}{1 - Pr\{X(\ell) = 0\}}.$$

The optimality equation for the entire horizon using the total reward criterion is as follows:

$$V_n(i) = \sum_{k=0,1,\dots,\bar{k}} \max_{a=0,1,\dots,k} \left\{ a \cdot R + U_n(i+a) \right\} dF(k), \text{ with} \quad (12)$$

$$U_n(i) = -E \left\{ c(i - X(\ell))^+ \right\} + E \left\{ V_{n-1}((i - X(\ell))^+) \right\}.$$

From (12), it is easily seen that if  $\Delta U_n(i) = U_n(i) - U_n(i+1)$  is non-decreasing in  $i$  for any  $n$ , or in other words if  $U_n(i)$  is concave, and non-increasing in  $i$ , then we have the following stationary policy for accepting customers:

$$a_n^*(i) = \arg \max_{a=0,\dots,M} \{\Delta U_n(i+a) \leq R\}.$$

**Theorem 4.2.**  *$V_n(i)$  and  $U_n(i)$  are concave and non-increasing in  $i$  for a fixed  $n$ , and the optimal acceptance policy for the Naive model is in the form of a critical level policy.*

From Theorem 4.2, there exists a critical accept-up-to policy ( $a_n^*$ ) that is independent of the number of customers in the system. This policy is very easy-to-use: accept  $\min\{a_n^* - i, k\}$  customers if there are  $i$  customers in the system.

#### 4.3.2.2 Leadtime Decisions of the Service-Sensitive Firm

In this section, we consider the case where the firm considers customers' sensitivity to the service level (i.e., the firm's past performance in meeting the quoted leadtimes). Similar to the case of the Naive firm, we formulate the problem as a MDP. However, now the state of the system is represented by  $(i, s)$ , where  $i$  is the number of orders in the system, and  $s$  is the service level at the beginning of the time period. Let  $V_n(i, s)$  be the maximal expected net benefit of the system over periods  $n$  ( $n \geq 0$ ) to 0 and  $V_0(i, s)$  be the expected reward to complete the service of the  $i$  customers remaining in the last period of the horizon, which is indicated as  $V_0(i, \cdot) = -[c \cdot i + E_{clearing}(i)]$  as before.

We model the problem using the total reward criterion, and the optimality equation is as follows:

$$V_n(i, s) = \sum_{k=0,1,\dots,\bar{k}} \max_{a=0,1,\dots,k} \{a \cdot R + U_n(i+a, s)\} dF_s(k), \quad \text{with} \quad (13)$$

$$U_n(i, s) = -E\{c(i - X(\ell))^+\} + E\{V_{n-1}(0, s^{up})\mathbf{I}_{(X(\ell) \geq i)} + V_{n-1}((i - X(\ell))^+, s^{down})\mathbf{I}_{(X(\ell) < i)}\},$$

where  $s^{up} = \min\{s+1, s_{max}\}$  and  $s^{down} = \max\{s-1, s_{min}\}$  are the system dynamics equations, and  $\mathbf{I}_{(\cdot)}$  is the indicator function.

Unlike the Naive firm, the order acceptance decisions of the Service-Sensitive firm depend on the service level. Let  $a_{n,s}^*(i)$  be the optimal accept-up-to level when there are  $i$  customers

in the system and the current service level is  $s$ , when there are  $n$  periods to go until the end of the horizon.

First, consider a given service level  $s$ . As the following example indicates, the optimal policy for the accept-up-to level is not stationary as the number of customers in the system changes.

**Example 4.1.** *Consider a five-period decision horizon where the arrivals follow a Poisson process and the maximum number of customers to arrive in a period is 10. We have 6 service levels  $S = \{0, 1, 2, 3, 4, 5\}$  with the corresponding arrival rates  $\{0, 5, 10, 15, 20, 25\}$ . The mean service rate is 5, and the lateness cost and revenue parameters are  $c = 1$  and  $R = 6$ , respectively. In this setting, the optimal accept-up-to decisions when the service level is fixed at 0 are given in Table 6.*

**Table 6:** Optimal accept-up-to level when the service level is 0

	Number of customers					
Time	0	1	2	3	4	5
1	10	10	10	10	10	10
2	10	10	10	10	10	10
3	5	4	3	2	10	10
4	5	4	3	2	1	0
5	5	4	3	2	1	0

In Table 6, the optimal order acceptance decisions in period 3 indicates an optimal accept-up-to level that neither increases nor decreases in the number of customers in the system.

Next, we look at how the optimal accept-up-to level changes as a function of the service level for a given number of customers in the system.

**Example 4.1.** *(Continued) The optimal accept-up-to decisions when there are six customers in the system prior to the accept/reject decisions are given in Table 7.*

As seen in Table 7, the optimal accept-up-to decisions are not necessarily monotonic (or convex or concave) in the service level.

Example 4.1 indicates that it is difficult to find a general structured order acceptance policy for the Service-Sensitive firm. Therefore, we consider a special case of this problem

**Table 7:** Optimal accept-up-to level when there are six customers in the system

	Service index					
Time	0	1	2	3	4	5
1	10	10	10	10	10	10
2	10	10	10	10	10	10
3	10	0	0	7	3	3
4	0	0	0	0	0	1
5	0	0	0	0	0	0

where the accepted orders are processed as batches and find that it has a nice structure, which can then suggest policies that might work well for the more general case. Batch processing would be applicable, for example, in chemical processing or the mass production of semiconductor chips with silicon wafers in furnaces, where the processing can be applied to many jobs without any negative effect on others. The optimal accept-up-to policy for batch processing is given in Theorem 4.3.

**Theorem 4.3.** *If service distribution is discretized uniform between  $[0, M]$  and no unethical customer acceptance policy is utilized (never accept to increase number of customers over  $M$ ), and the firm’s service performance is modelled by 3 service levels, the optimal order acceptance policy  $(a_n^*(s))$  in each period is a threshold policy, or an “accept-up-to” policy, that depends on the current service level. When there are  $i$  customers in the system and the service level is  $s$ , the optimal number of accepted customers is  $(a_n^*(s) - i)^+$ . The optimal accept-up-to policy has a convex structure in the service level.*

This policy is more complex than in the case without incorporating service, since it is non-monotonic in the service level. However, for this special case, it helps to reduce the search space for the best policy, and it also suggests a structure for a policy that might work well in practice. We use the structured policy given in Theorem 4.3 as one heuristic for the general Service-Sensitive model. We examine the performance of this heuristic in Section 4.4.

#### 4.3.2.3 Heuristics

Since the general Service-Sensitive model with finite capacity has a complex optimal policy, it may be difficult to solve for large instances. Thus, we have developed several heuristics for the general problem that are more computationally efficient than solving to optimality. The heuristics are briefly described below.

- *Myopic Heuristic:* The myopic heuristic incorporates service but only looks ahead a small number of periods instead of solving the entire dynamic program optimally. We implement and test this heuristic for varying number of periods in numerical experiments and find the critical accept-up-to level for each service level.
- *Marginal Cost/Benefit (MCB) Heuristic:* In this heuristic we examine accepting one more customer, where the estimated value accounts for the immediate revenue increase but also the expected increase in penalty cost and impact of a potential service decrease on the future. Given that  $i - 1$  customers are already accepted, we estimate the marginal value of accepting the  $i^{\text{th}}$  customer at the current service index  $s$  by:

$$\begin{aligned} \Delta\Pi(i) = & -Pr\{X(\ell) < i - 1\} \cdot c + Pr\{X(\ell) > i - 1\} \cdot R \\ & + Pr\{X(\ell) = i - 1\}(R - c - 2\alpha R\beta(s)) \end{aligned}$$

The first and second terms of  $\Delta\Pi(i)$  capture the case when the acceptance of the  $i^{\text{th}}$  customer does not change the future service level, and an extra cost or revenue is obtained due to accepting this customer. The third term corresponds to the case when the acceptance of customer  $i$  results in a drop in service, which causes the arrival rate of future customers to decrease. Therefore the firm loses future expected revenue of  $2\alpha R\beta(s)$  in the next period, where  $\beta(s)$  is a parameter indicating the proportion of lost revenue for each service level, therefore  $\beta \in [0, 1]$ . (If the service level is  $s_{min}$  or  $s_{max}$  then we replace  $2\alpha$  with  $\alpha$  because for those service levels there is a maximum drop of one level.) Having  $\beta(s)$  depend on  $s$  ensures that the policies of the heuristic depend on  $s$  just as the optimal policies do; it is also reasonable to expect that the risk of losing future arrivals may be different for different service levels. We choose



an increment and search over the best discrete values of  $\beta$  for each service level. The set of  $\beta$  values that provide the highest expected profit is used to find accept-up-to levels. If there are too many service levels, it may be time-consuming to search for  $\beta$  for all of them. In that case, it is possible to reduce the search time by assuming that  $\beta$  is the same for some service levels. For instance, it may be reasonable to assume that for any  $s$  when  $\lambda(s) > \mu + 2\sigma$ , the same  $\beta$  may apply.

- *Convex Accept Levels (CAL) Heuristic*: Since we know that the special case of batch processing results in optimal decisions that are convex in the service levels, we search specifically for convex policies for the general problem. While we do not have an efficient way of identifying the very best convex policy, we use the following heuristic for finding a high quality convex policy: Search for the best accept level for the first service level and reduce the search afterwards to be convex around that initial value and find the accept-up-to levels by this procedure for each period.

## 4.4 Numerical Experiments

### 4.4.1 Experiment Details

We perform numerical experiments to gain additional insights on the impact of the customers' sensitivity to the service performance. Recall that the Service-Sensitive firm uses a policy that is optimal considering the customers' sensitivity and obtains a profit of  $\Pi^*(L)$ , whereas the Naive firm uses a policy that is optimal for a constant customer arrival rate and obtains a profit of  $\Pi(L)$ . As a performance metric, we look at the percentage improvement in profit from using the Service-Sensitive versus the Naive model, i.e.  $\frac{\Pi^*(L) - \Pi(L)}{\Pi(L)} \times 100$ .

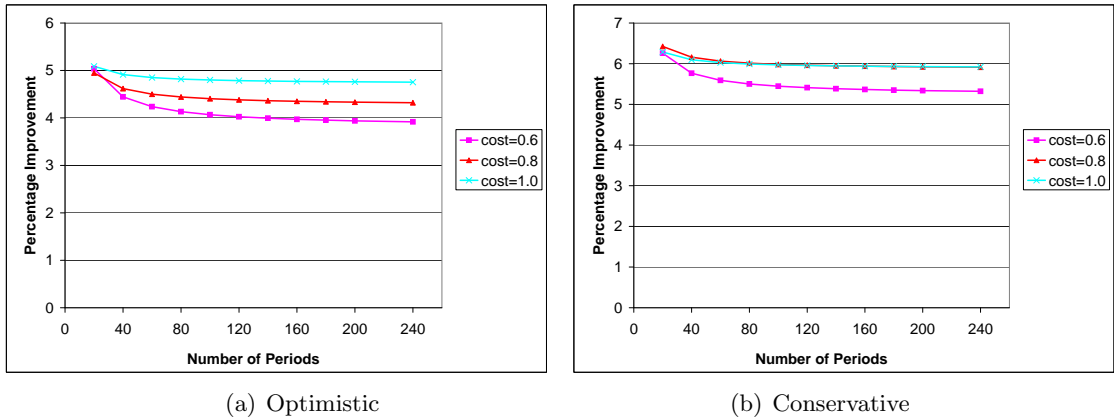
We assume that both the arrival and the service processes are Poisson. The arrival rate of the customers changes with the service level according to  $\lambda(s) = \lambda_0 + \alpha \cdot s$ , where  $\lambda_0$  is the base arrival rate and  $\alpha$  is the sensitivity of the customers to the service performance of the firm. We consider two cases for the arrival rate, which the Naive firm believes to be constant: *conservative* and *optimistic*, where the Naive firm believes that the constant arrival rate is equal to the mean and the maximum arrival rate, respectively. We assume that at most 10 customers arrive at the beginning of each period to keep the problem size

manageable while finding the optimal policy by dynamic programming techniques.

#### 4.4.2 Results

Many companies use score cards with five levels (e.g., ranging from 0=below expectations to 4=well above expectations) to evaluate their suppliers (see, for example Institute of Chartered Accountants [75] and Samtec Inc., Supplier Quality Assurance Manual [78]). Hence, for most of the experiments, we model the firm's service performance by 5 service levels, 0 indicating the lowest and 4 indicating the highest performance of the firm. Revenue from a customer is 6 for all computations, and until otherwise stated we assume that the mean service rate is 5 and the mean arrival rate changes from 1 to 9, which implies  $\lambda_0 = 1$  and  $\alpha = 2$  if there are 5 service levels.

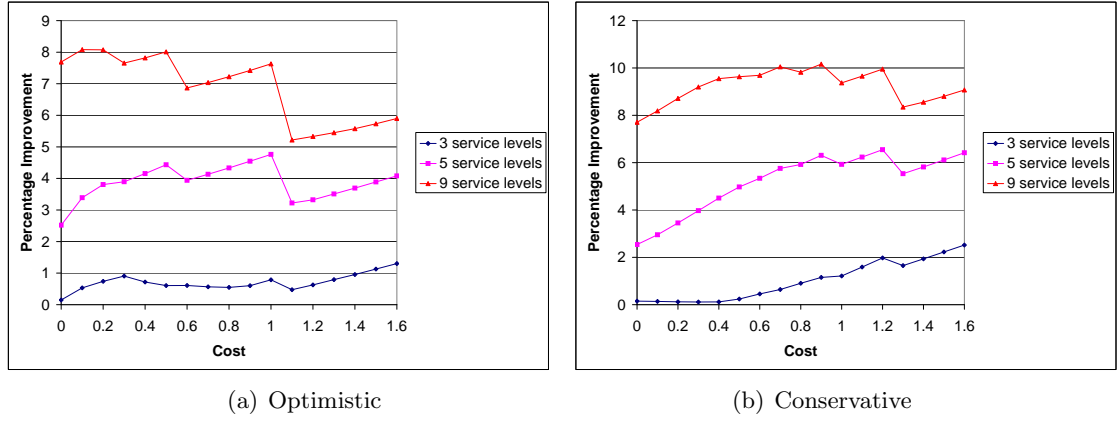
Figure 9, shows how the percentage improvement in profit changes as the horizon length increases when we solve both the Service-Sensitive and the Naive models optimally. At each cost level, the percentage improvement decreases quickly in the first part of the horizon and then stabilizes. The decrease at the beginning is largely due to the better performance of the Service-Sensitive model in the vicinity of the terminal period. As seen in Figure 9, the initialization effect becomes insignificant after 200 periods, therefore, we use 200 periods in the remainder of our experiments.



**Figure 9:** Percentage Improvement of Optimal over Naive with the Horizon Length

In Figure 10, we examine how the percentage improvement in profits from considering service changes as the number of service levels increases. The percentage improvement (i.e., the benefit of considering the customers' service sensitivity) is higher when there are

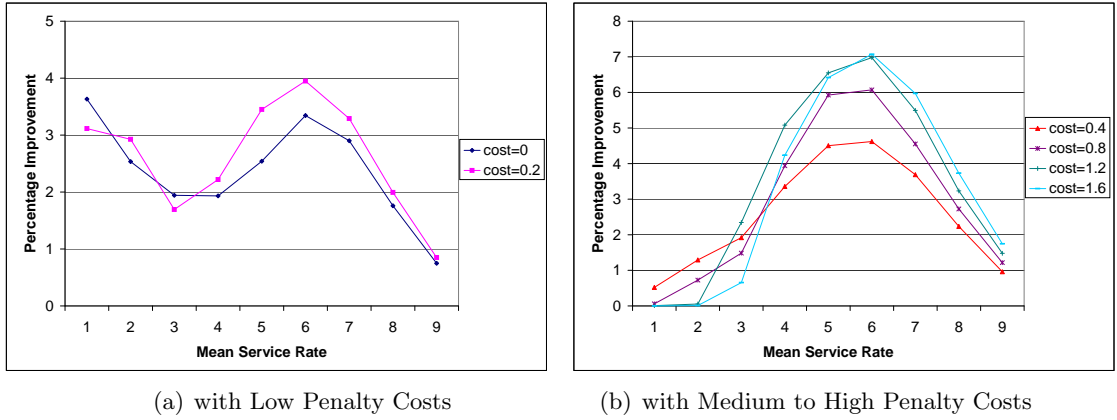
more service levels both for the conservative and the optimistic cases. The percentage improvement is as high as 10% (8%) for the conservative (optimistic) case, which can be explained by the fact that higher number of service levels allows the policy to control the arrival rate of the customers more closely.



**Figure 10:** Percentage Improvement of Optimal over Naive with the Number of Service Levels

Next, we look at the impact of capacity (i.e., mean service rate) on the percentage improvement in profits. When the penalty cost is medium to high (Figure 11(b)), the percentage improvement follows a concave structure, i.e., the improvement is small for tight or abundant capacity, and large for medium capacity. Intuitively, when the capacity is abundant, a high number of customers may be accepted and served on time, and there is less difference between the decisions (and the profits) of the Naive and the Service-Sensitive firms. Similarly, when the capacity is tight and the cost is not very small, both the Naive and the Service-Sensitive firms accept fewer customers to avoid delay penalties. However, when the capacity is medium, there is more room for error; in particular, by accepting a higher number of customers than the Service-Sensitive firm, the Naive firm loses a significant amount of future business. When the penalty cost is small (Figure 11(a)), we see that the improvement first decreases in the mean service rate, and then follows a concave structure. When the cost is small, the Naive firm tends to accept a significantly higher number of customers, even for low capacity levels, since the revenues outweigh the penalty costs. This leads to a significant loss of future business, resulting in a high difference between Service-Sensitive and Naive profits. As the service rate slightly increases, the loss of future business

decreases and hence, the profit difference also decreases. As the service rate increases further, we see a similar concave structure in the percentage improvement of profits as in the medium-high cost case, and similar explanations hold.



**Figure 11:** Percentage Improvement of Optimal over Naive with Service rate

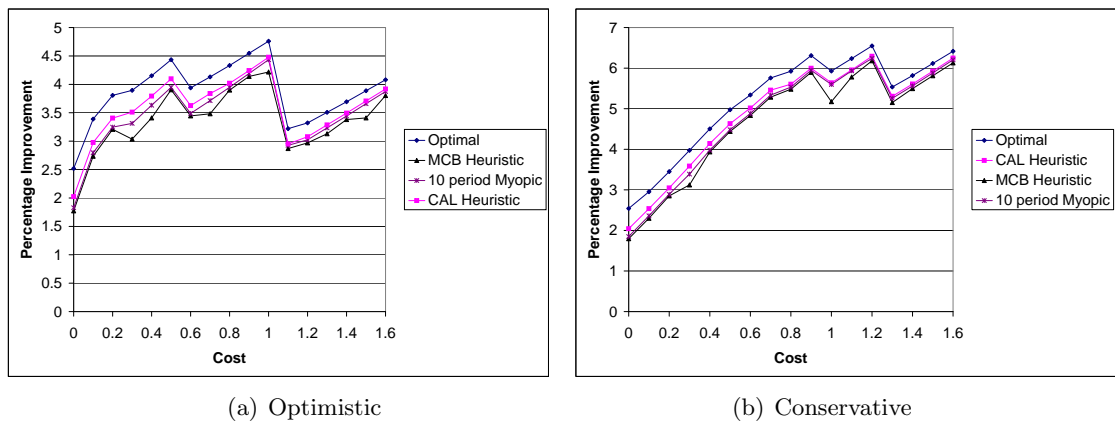
In the results this far, we calculate the optimal policies and profits by dynamic programming. As the number of periods increases, the solution time increases rapidly due to the size of the problem. For large problems, efficient heuristics are needed to incorporate the service performance into the leadtime and order acceptance decisions. We compare the optimal policy and the heuristic policies that are mentioned in Section 4.3.2.3 to the Naive case, for 200 periods and 5 service levels.

For the *Myopic* heuristic, we tried several horizon lengths (namely, 5, 10, and 15). We found the marginal improvement of using 15 periods over 10 periods to be very small, therefore we use 10 periods in our experiments. In our implementation, we determine the best accept-up-to policy using the last periods of the dynamic program and apply this stationary policy to all periods. In the *Marginal Cost/Benefit* heuristic, we consider a constant  $\beta(s)$  if the arrival rate at service index  $s$  is within some predefined range (e.g.  $(\cdot, \mu - \sigma)$ ,  $[\mu - \sigma, \mu + \sigma]$ ,  $(\mu + \sigma, \cdot)$ ) of the mean service rate  $\mu$ . This approach reduces the number of parameters to be considered (computed) in the heuristic. For the case of 5 service levels, we have three  $\beta(s)$  values which we estimate by running 100 repetitions of a 200-period simulation over the all possible values of  $\beta$  using an 0.1 increment and by looking at the average profit values. For the 200-period problem, the *Marginal Cost/Benefit*

heuristic takes 10% of the time that the optimal dynamic program requires to find the optimal decisions, and note that as the horizon gets longer the solution time of the heuristic increases linearly whereas the solution time of the dynamic program increases exponentially.

In the *Convex Accept Levels* heuristic, the accept-up-to levels are found by the same dynamic programming technique used to find the optimal policies for the Service-Sensitive and Naive cases, only this time the policy is forced to be a threshold level which has a convex structure in the service index. Although there is some improvement in the computation time, the requirement of solving the entire dynamic program is still a concern in terms of run time.

In general, we find that the Convex Accept Levels heuristic performs best, closely followed by the Myopic and the Cost/Benefit Heuristics (they result in 97%, 96%, and 95% of the optimal profits, respectively, for medium-high cost). The performance of the *Marginal Cost/Benefit* Heuristic is slightly worse in our experiments, but its performance may be improved by a better estimation of the  $\beta$  values, for instance, by using a different  $\beta$  for every service level. It is also very encouraging that the behavior of the optimal and the heuristic policies are very similar, which suggests that the heuristics capture the essence of the optimal policy.



**Figure 12:** Percentage Improvements of Optimal and Heuristic over Naive with Increasing Cost

## 4.5 Conclusions

Recent examples of several companies emphasize the importance of leadtime performance on the viability of the firm. The outcomes from ignoring service when quoting leadtimes or accepting customers has ranged from the loss of millions of dollars paid in penalty costs, the loss of customers due to the decrease in the reputation of the firm, and even to firms going out of business (e.g., [43], [50] and eToys). Clearly when managing their customer demand, a firm ignores service considerations at their peril.

In this chapter we consider leadtime and order acceptance decisions in a manufacturing firm with stochastic arrivals and production processes, where a manufacturer considers past service when making decisions. We assume that customers are aware of past service (such as through internet tracking sites or by reputation), that a customer's probability of placing an order decreases with increasing leadtimes, and that past service determines the level of future customer arrivals. A Service-Sensitive firm considers past performance in their demand management policies, while a Naive firm ignores past performance in their policies; both firms consider expected revenue as well as expected penalty costs due to late orders.

When there is infinite production capacity, the long-run service is the proportion of orders completed on-time, and we assume that the processing times are exponential. This allows us to find a closed-form expression for leadtime quotation. We show that the optimal leadtime to quote that accounts for past performance is

- more conservative (i.e., longer) than the optimal leadtime that ignores it, and
- always positive, which means that a Service-Sensitive firm would never quote unethical leadtimes.

This last result is important since it is possible that a Naive firm will quote unethical leadtimes (i.e., leadtimes of zero) when the revenues are sufficiently high.

We also study demand management decisions when capacity is limited. In this case, we assume that the leadtime is an industry-dictated standard, therefore considering whether or not to accept a customer is equivalent to quoting a leadtime equal to the standard leadtime. We also assume that decisions are made at discrete time intervals (e.g., due to high review

costs), service is measured with discrete levels, and the proxy for service is whether all orders were completed on-time in a given period.

For the general finite capacity model with a Service-Sensitive firm, we demonstrate that there is not a structured acceptance policy with the number of customers or service levels, even with Poisson arrivals and exponential service times. However, for the special case where production is done in batches, we show that the optimal acceptance policy is of a threshold type with the number of customers, and this policy is further convex in the service levels. (The structure of this result also helps to inspire a heuristic for the more general problem.) For the corresponding Naive model, we prove that a threshold acceptance policy in the number of customers is optimal.

We develop several heuristics to solve the general Service-Sensitive problem when it is computationally expensive to find the optimal solution using dynamic programming. The heuristics perform within 95-97% of the optimal solution for medium-high lateness costs.

In our numerical experiments, we find that considering service where there is limited capacity can have an impact of more than 2 – 6% in the profit over the Naive case. The benefit from incorporating service is high when demand is close to the capacity, and having a higher number of service levels tends to further increase the benefit compared to the Naive case.

This chapter is the first work in the production and leadtime literature to incorporate past performance on customers' decisions; it has shown that service matters significantly in leadtime quotation and demand management, and it is a starting point for many more applications that can be considered in this area. Better service leads to higher customer satisfaction, increases brand loyalty, and leads to higher profits for the firm.

## CHAPTER V

### DYNAMIC SWITCHING TIMES FOR SEASON AND SINGLE TICKETS IN SPORTS AND ENTERTAINMENT

#### *5.1 Introduction*

Revenue Management (RM) has made great strides in improving the bottom line of many firms, especially in airlines (Smith et al. [70]), hotels (Lieberman [52]), and rental car agencies (Geraghty and Johnson [37]), where RM is recognized as a key factor in the firm's viability and success. In these industries, there is often a limited or fixed capacity, and firms are able to segment the market according to differing customer needs for particular products or services. Revenue Management is a set of tools to help mathematically determine decisions such as the right prices or inventory to make available so as to maximize profit.

However, there are many other industries that offer a rich set of RM-type problems that have not been fully addressed. One of these is the sports and entertainment (S&E) industry, where tickets are sold in advance to an event at a venue such as a sports stadium or theater. Like the airlines, the capacity for an event is generally fixed in advance, there are high fixed costs to operate the venue and low marginal costs to selling additional tickets, and the market can be segmented into different kinds of customers.

In S&E, one important segmentation of the market is that some customers buy season packages, or bundles of tickets to events during the season, while others buy individual tickets to performances. Season ticket holders are important to the success of the organization, since they are more likely to donate to the organization, buy apparel, or renew tickets in the future. They are also desirable customers since they commit to a bundle of tickets in advance, which can offer greater cash flow to an organization or commitments upon which to base future operational decisions. Most S&E firms offer season tickets first, and they open purchasing for single tickets at a later date but before the start of the season. A basic trade-off is that the firms want to capture as much of the demand for bundles of ticket,



while still allowing enough time for individual ticket purchases when bundle demand will not be sufficient to sell out the stadium, as it usually is not.

There are many interesting questions in Revenue Management in S&E. In this chapter, we study the specific question of timing the switch from selling bundles to selling single tickets. This is a problem motivated by our discussions with several large S&E firms, where the timing of the sales, or in some cases, the timing of the *promotion* of sales to the public is of key interest. The problem we study also is relevant to other industries where revenue management applies. For instance, many hotels make accommodations for group bookings for weddings or conferences, but the commitments must be made in advance and the unsold rooms are released to the general public in advance of the travel date for bookings by individuals. It is also possible to sell bundled capacity and smaller units of capacity in industries such as manufacturing, where contracts may be negotiated with prioritized clients for larger volumes of capacity.

A key aspect to the problem we study is that the bundled and single tickets are sharing the same, limited capacity. In addition, after sales of single tickets are allowed, there may be multiple events for sale simultaneously. These characteristics, along with the desire to *dynamically* determine the timing decision when demand is stochastic, necessitated the development of new models for the RM decision-making. Although the mathematics are complex, we find that the structure of the problem leads to an optimal timing policy that is relatively easy to understand and implement. The resulting policy defines a set of threshold pairs of times and remaining inventory which determine the switch from bundles to single ticket sales. After each bundle sale, if the current time is less than the corresponding threshold, then the switch is made to selling individual tickets. We describe an algorithm that will compute the threshold pairs, and we demonstrate the value of the dynamic timing decision. We are able to generalize our results in several ways, including allowing the demand rates for the bundles and single-tickets to depend on time.

In the next section we describe the relevant literature and identify our contribution. In Sections 5.3 and 5.4, we introduce the assumptions and the model, and present key results for the base case. We generalize the model in Section 5.5, and we demonstrate some

numerical examples in Section 5.6. Finally, we conclude in Section 5.7 and offer several directions for promising research in RM in S&E.

## ***5.2 Literature Review***

In the airlines, revenue management research include how to determine the overbooking levels for each fare-class (Littlewood [54] and Belobaba [6, 7]) and the bid-prices for each leg of a network (Williamson [88] and Talluri and van Ryzin [79]); recent applications in airline RM include Bertsimas and Popescu [8] and Karaesmen and Ryzin [45]. Unlike S&E, the airline industry has a network structure where demand is for an origin and destination pair, which may include multiple choices of paths for the consumer. When group purchases are considered in the airlines, they are primarily for groups of individuals purchasing tickets on one plane, rather than a single individual purchasing multiple tickets over time, and there is very limited literature on group sales as stated in Yuen [89] and Farley [26]. The main focus in airline RM is on determining prices or seat allocations, possibly across multiple segments of customers, where customers purchase one ticket, rather than the timing of bundles and single sales. The most similar sale to season tickets is the offering of “flexible products”, where a single individual buys the option of two or more flights at a time and assigned to one of them later by the carrier (Gallego and Phillips [34]).

There have been several papers in RM of airline and retail industries that focus on pricing as a function of time. Gallego and van Ryzin [35] study pricing of a set stock of products to be sold by a deadline and use intensity control to identify optimal prices as a function of the stock level and remaining time. They also show the asymptotically optimality of the policies with at most one price change as the volume of sales increase. In the S&E industry, most organizations keep prices as announced throughout the selling period, which is known as price stickiness in the entertainment industry (see Courty [16]). Therefore pricing that are used in retailing are less applicable to the entertainment industry. In the S&E, timing of different kinds of products is more common than timing of a price change.

A closely related paper to our work is Feng and Gallego [29], which determines the

optimal dynamic time to switch from one predetermined price to a second higher or lower predetermined price so as to maximize revenue by selling a given stock over a finite time horizon. Demand is assumed to be stochastic and demand rate is higher for the lower price, and the optimal timing policy is shown to be a time threshold depending on the remaining stock amount. The restrictions of one price change and time-invariant demand intensities are relaxed in Feng and Gallego [28], and an efficient algorithm to find the optimal value functions and the optimal pricing policy is provided. Like the latter papers, we use predetermined prices, but a main difference in our work is that we focus on switching sales from bundles of tickets to single tickets. A second important factor in comparison to [29] and follow-on papers is that when we switch to selling singles, the bundles split into multiple simultaneous processes.

An initial version of the switching problem between bundles and singles is studied in S&E in Drake et al. [22]. However, in that paper, they specifically focus on a static timing decision, as is done in some organizations, where the switching date to single tickets is announced in advance to the public. In this work, we study the *dynamic* switching time, where the time may be determined by the sales-to-date. Although this complicates the mathematics, it is important, since some organizations dynamically select their switching or promotions times based on past demand. In [22], they assumed a linear Markovian death process, but in this chapter we generalize the demand function to be any Poisson process, so the techniques for analysis are quite different.

It is also important to point out that there has been analysis related to improving revenue in S&E in other disciplines. A number of papers have looked at pricing decisions within venue but did not consider bundling. For example, Leslie [51] and Rosen and Rosenfield [66] studied revenue-maximizing ticket prices for different seat qualities but neither of these studies considers the bundling of tickets. The most relevant work in the economics and marketing literature that considers the selling of bundled commodities are: Venkatesh and Mahajan [85], Venkatesh and Kamakura [84], McAfee et al. [56], Salinger [67] and Bakos and Brynjolfsson [4], but these papers focus on the pricing of the bundles, not the timing of decisions.

### 5.3 Assumptions and Notation

Let  $M \in \mathbb{Z}^+$  be the number of seats available for sale for each performance, and  $T \in \mathbb{R}^+$  be the selling period. In the S&E data we have seen, season tickets are rarely bought after the season begins, and the switch to selling singles is also made before the season starts in every organization with whom we have worked. Thus, we focus on the selling horizon before the season begins and assume that the selling period ends when the first performance takes place. The selling period begins with first offering tickets as a bundle at price  $p_B$  and then switching to selling performance tickets individually at  $p_i$  for performance  $i$ , for  $i = 1, 2$ . We assume that these prices are predetermined at the beginning of the selling season, which is true for most organizations, especially during the time preceding the start of the season. Note organizations may have multiple classes but here we focus on the two different products, with average price for each.

We assume that market segments (bundles and singles) are independent. This is supported by discussions with professional sports teams (Depaoli [20]), and it is also a common assumption for many models in revenue management. We initially assume constant demand rates with time for each of the bundled and single-ticket processes. In the second part of the chapter, we extend the model and results to allow demand to depend on time. In the extensions, the rates can also be used to proxy substitution among segments, by allowing the demand rate for singles to be higher earlier in the season.

We assume that for each price, there is a corresponding Poisson process:  $N_B(s)$ ,  $0 \leq s \leq t$ , with known constant intensity  $\lambda_B$  for the bundled performances;  $N_1(s)$ ,  $0 \leq s \leq t$ , with known constant intensity  $\lambda_1$ , and  $N_2(s)$ ,  $0 \leq s \leq t$ , with known constant intensity  $\lambda_2$  for the two single performances, respectively. The state of the system is indicated by the elapsed time  $t$  and the remaining inventory level at time  $t$ ,  $n(t)$ .

We define  $r_B = \lambda_B p_B$  and  $r_i = \lambda_i p_i$  as the revenue rate from the bundled and individual ticket sales of the single performances  $i = 1, 2$ , respectively. We assume that the expected revenue rate for the bundle is higher than the sum of the expected revenue rates of the single tickets, i.e.,  $r_B > r_1 + r_2$ . Otherwise, switching immediately would be optimal for all states, and it is not relevant to study the optimal time to switch. This assumption can also be

intuitively validated by the fact that the revenue for each bundle sale can include intangibles such as donations to the organization (sometimes required for season ticket purchases) or the value of early commitment and guaranteed revenues.

#### 5.4 Model and Results

The expected revenue over  $[t, T]$  is given by two expressions:  $\Pi(t, n(t))$  and  $V(t, n(t))$ .  $\Pi(t, n(t))$  is the expected revenue when  $n(t)$  items are available for sale over  $[t, T]$  while tickets are sold individually (which means that the switch from bundles has already occurred) and it is given by:

$$\Pi(t, n(t)) = p_1 E[(N_1(T) - N_1(t)) \wedge n(t)] + p_2 E[(N_2(T) - N_2(t)) \wedge n(t)],$$

where  $(x \wedge y)$  indicates the minimum of the two terms.  $V(t, n(t))$  is the *optimal* expected revenue over  $[t, T]$ , when  $n(t)$  items are available for sale over  $[t, T]$  and  $\tau$  is the best switching time to individual ticket sales. It is represented by:

$$V(t, n(t)) = \sup_{\tau \in \mathcal{T}} E[p_B((N_B(\tau) - N_B(t)) \wedge n(t)) + \Pi(\tau, n(\tau))],$$

where  $\mathcal{T}$  is the set of switching times  $\tau$  satisfying  $t \leq \tau \leq T$  and  $n(\tau) = [n(t) - N_B(\tau) + N_B(t)]^+$ , and  $x^+ = \max\{0, x\}$ .

At time  $t$ , if we can compare the expected revenue of switching immediately from selling bundles to the expected revenue of delaying the switch to a time  $\tau$  ( $t \leq \tau \leq T$ ), then we can decide whether delaying the switch further than time  $t$  is beneficial or not. At the state  $(t, n(t))$  the expected revenue of switching at  $t$  is:  $\Pi(t, n(t))$ , and the expected revenue of switching at  $\tau$  for  $\tau \geq t$  is  $E[p_B((N_B(\tau) - N_B(t)) \wedge n(t)) + \Pi(\tau, n(\tau))]$ .

To compare these two expected values, we need a tool to measure the infinitesimal effect of the delay. Let us define the infinitesimal generator  $\mathcal{G}$  with respect to the Poisson process

$(t, N_B(t))$  for a uniformly bounded function  $g(t, n)$  as:

$$\begin{aligned}
\mathcal{G}g(t, n) &= \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} E[g(t + \Delta t, n - N_B(\Delta t)) - g(t, n)] \\
&= \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} \sum_{k=0}^{\infty} [g(t + \Delta t, (n - k)^+) - g(t, n)] \frac{(\lambda_B \Delta t)^k}{k!} e^{-\lambda_B \Delta t} \\
&= \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} [(g(t + \Delta t, n) - g(t, n))(1 - \lambda_B \Delta t) + (g(t + \Delta t, n - 1) - g(t, n))\lambda_B \Delta t] \\
&= \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} (g(t + \Delta t, n) - g(t, n)) + \lim_{\Delta t \rightarrow 0} \lambda_B (g(t + \Delta t, n - 1) - g(t + \Delta t, n)) \\
&= \frac{\partial g(t, n)}{\partial t} + \lambda_B [g(t, n - 1) - g(t, n)].
\end{aligned}$$

Applying  $\mathcal{G}$  to the function  $\Pi(t, n(t))$  gives the immediate loss of single ticket revenue from delaying the switch from selling bundles to selling singles. Specifically,  $\mathcal{G}\Pi(t, n(t)) = \frac{\partial \Pi(t, n(t))}{\partial t} + \lambda_B [\Pi(t, n(t) - 1) - \Pi(t, n(t))]$ , which is composed of two parts:  $\frac{\partial \Pi(t, n(t))}{\partial t}$ , which is the loss of revenue due to elapsed time, and  $\lambda_B [\Pi(t, n(t) - 1) - \Pi(t, n(t))]$ , which is the loss of revenue due to the decrease in inventory to be sold as singles. But during the time when the switching is delayed, the Poisson process for bundles  $(t, N_B(t))$  is active, and it generates revenue at the rate  $\mathcal{G}E[p_B((N_B(\tau) - N_B(t)) \wedge n(t))] = \lambda_B p_B$ . Therefore, the net marginal gain (or loss) for delaying the switch from bundles to singles at state  $(t, n(t))$  is given by:

$$\mathcal{G}\Pi(t, n(t)) + \lambda_B p_B = \frac{\partial \Pi(t, n(t))}{\partial t} + \lambda_B [\Pi(t, n(t) - 1) - \Pi(t, n(t))] + \lambda_B p_B.$$

By Dynkin's Lemma (Rogers and Williams [65]), we have the following two martingales for any  $s \geq t$ :

$$\Pi(s, n(s)) - \Pi(t, n(t)) - \int_t^s \mathcal{G}\Pi(u, n(u)) du, \quad (14)$$

$$p_B((N_B(s) - N_B(t)) \wedge n(t)) - p_B((N_B(t) - N_B(t)) \wedge n(t)) - \int_t^s \lambda_B p_B \mathbf{1}_{\{n(u) > 0\}} du, \quad (15)$$

where  $\mathbf{1}_{\{n(u) > 0\}}$  is an indicator function. Since, the expected value of these martingales at any time  $s$  is equal to their expected value at the starting time  $t$ , we have:

$$\Pi(s, n(s)) - \Pi(t, n(t)) = E \int_t^s \mathcal{G}\Pi(u, n(u)) du, \quad (16)$$

$$E[p_B((N_B(s) - N_B(t)) \wedge n(t))] = E \int_t^s \lambda_B p_B \mathbf{1}_{\{n(u) > 0\}} du. \quad (17)$$

By the optional sampling theorem (Karatzas and Shreve [46]), we can replace  $s$  in (16) and (17) with any stopping time  $\tau \geq t$ . Therefore, adding equations (16) and (17) for a stopping time  $\tau$ , we get:

$$\begin{aligned} E[p_B((N_B(\tau) - N_B(t)) \wedge n(t)) + \Pi(\tau, n(\tau))] - \Pi(t, n(t)) \\ = E \int_t^\tau [\mathcal{G}\Pi(u, n(u)) + \lambda_{BPB} \mathbf{1}_{\{n(u) > 0\}}] du. \end{aligned} \quad (18)$$

Note that the left-hand side of (18) is the expected revenue gained over  $\Pi(t, n(t))$  by delaying the switch from  $t$  to  $\tau$ , and we can quantify it by using  $\mathcal{G}$ , as shown in the right-hand side. Therefore, we know that delaying the switch to  $\tau$  from  $t$  is beneficial if  $E \int_t^\tau [\mathcal{G}\Pi(u, n(u)) + \lambda_{BPB} \mathbf{1}_{\{n(u) > 0\}}] du > 0$ .

Taking the supremum of both sides in (18) over all stopping times  $t \leq \tau \leq T$ , and defining:

$$\tilde{V}(t, n(t)) = \sup_{t \leq \tau \leq T} E \int_t^\tau [\mathcal{G}\Pi(u, n(u)) + \lambda_{BPB} \mathbf{1}_{\{n(u) > 0\}}] du, \quad (19)$$

we get that  $V(t, n(t)) = \Pi(t, n(t)) + \tilde{V}(t, n(t))$ . This implies that the optimal revenue over  $[t, T]$  consists of two parts: the revenue from the immediate switch (selling single tickets until the end of horizon) and the additional revenue from delaying the switch further in time. Since  $\tilde{V}(t, n(t))$  is also given by:

$$\tilde{V}(t, n(t)) = \sup_{t \leq \tau \leq T} E[p_B((N_B(\tau) - N_B(t)) \wedge n(t)) + \Pi(\tau, n(\tau))] - \Pi(t, n(t)), \quad (20)$$

it is obvious that  $\tilde{V}(t, n(t)) \geq 0$  for any  $0 \leq t \leq T$  and  $0 \leq n(t) \leq M$ . In particular,  $\tilde{V}(t, 0) = 0$  for all  $0 \leq t \leq T$  and  $\tilde{V}(T, n(t)) = 0$  for all  $0 \leq n(t) \leq M$ . Moreover, equation (20) indicates that when  $\tilde{V}(t, n(t)) = 0$ , delaying the switch further is not optimal, whereas  $\tilde{V}(t, n(t)) > 0$  implies a revenue potential from delaying the switch.

To compute  $\tilde{V}(t, n(t))$ , we introduce a function  $\bar{V}(t, n(t))$ , which can be derived recursively, and is identical to  $\tilde{V}(t, n(t))$  when a number of conditions are satisfied. Obviously,  $\bar{V}(T, n(t)) = 0$  and  $\bar{V}(t, 0) = 0$  must be in the list of conditions. Also, since  $\tilde{V}(t, n(t))$  determines whether it is optimal to switch immediately or not,  $\bar{V}(t, n(t))$  must also imply the switching decision. Formally,

**Theorem 5.1.** *Suppose there exists a function  $\bar{V}(t, n(t))$  such that  $\bar{V}(t, n(t))$  is continuous and differentiable with right continuous derivatives in  $[0, T]$  for each fixed  $n(t)$ . In addition, if  $\bar{V}(t, n(t))$  satisfies:*

$$(i) \bar{V}(t, n(t)) \geq 0, \quad 0 \leq t \leq T \text{ and } 0 \leq n(t) \leq M;$$

$$(ii) \bar{V}(T, n(t)) = 0 \text{ for } 0 \leq n(t) \leq M \text{ and } \bar{V}(t, 0) = 0 \text{ for } 0 \leq t \leq T;$$

$$(iii) \bar{V}(t, n(t)) = 0 \Rightarrow \mathcal{G}(\bar{V} + \Pi)(t, n(t)) + \lambda_B p_B \leq 0, \quad 0 \leq t \leq T \text{ and } 0 \leq n(t) \leq M;$$

$$(iv) \bar{V}(t, n(t)) > 0 \Rightarrow \mathcal{G}(\bar{V} + \Pi)(t, n(t)) + \lambda_B p_B = 0, \quad 0 \leq t \leq T \text{ and } 0 \leq n(t) \leq M;$$

then  $\bar{V}(t, n(t)) = \tilde{V}(t, n(t))$ .

The proofs for Theorems 5.1 and 5.2 are essentially that of Feng and Xiao [30] and they are provided in Appendix.  $\bar{V}(t, n(t))$  enables us to decide whether to delay the switch further than  $t$  is beneficial or not. The net marginal gain from delaying,  $\mathcal{G}\Pi(t, n(t)) + \lambda_B p_B$ , is the main term that defines the behavior of  $\bar{V}(t, n(t))$ , and this term will be addressed more closely in the following lemma. First, noting that  $E[(N_i(T) - N_i(t)) \wedge n(t)] = \sum_{k=1}^{n(t)} P[N_i(T) - N_i(t) \geq k]$ , we can express  $\Pi(t, n(t))$  for  $n(t) \geq 1$  as:

$$\Pi(t, n(t)) = p_1 \sum_{k=1}^{n(t)} P[N_1(T) - N_1(t) \geq k] + p_2 \sum_{k=1}^{n(t)} P[N_2(T) - N_2(t) \geq k]. \quad (21)$$

**Lemma 5.1.** *The net marginal gain from delaying for  $0 \leq t \leq T$  can be written as:*

$$\begin{aligned} \mathcal{G}\Pi(t, n(t)) + \lambda_B p_B &= (r_B - r_1 - r_2) + p_1(\lambda_1 - \lambda_B)P[N_1(T) - N_1(t) \geq n(t)] \\ &\quad + p_2(\lambda_2 - \lambda_B)P[N_2(T) - N_2(t) \geq n(t)]. \end{aligned}$$

See Appendix for details of the proof. Since  $\lambda_B > \lambda_i$  for  $i = 1, 2$ , clearly  $\mathcal{G}\Pi(t, n(t)) + \lambda_B p_B$  is increasing in  $t$  and  $n$ . Noting that  $\mathcal{G}\Pi(T, n(T)) + \lambda_B p_B = r_B - r_1 - r_2$  when  $n(T) \geq 1$ , we can conclude that  $\mathcal{G}\Pi(t, n(t)) + \lambda_B p_B \leq r_B - r_1 - r_2$ .

If  $r_B \leq r_1 + r_2$ , we have  $\mathcal{G}\Pi(t, n(t)) + \lambda_B p_B \leq 0$  for all  $(t, n(t))$  which implies  $\bar{V}(t, n(t)) = 0$  for all  $(t, n(t))$  too, because it satisfies the conditions at Theorem 5.1 ( $\mathcal{G}\Pi(t, n(t)) + \lambda_B p_B = \mathcal{G}(0 + \Pi(t, n(t)) + \lambda_B p_B \leq 0$ ). This result validates the assumption that  $r_B > r_1 + r_2$  is required for bundle sale option to be considered.



Although we have demonstrated the existence of the alternate function  $\bar{V}(t, n(t))$ , the issue of how to calculate it for any  $(t, n(t))$  pairs still remains. From condition (iv), we know that  $\mathcal{G}\bar{V}(t, n(t)) = -\mathcal{G}\Pi(t, n(t)) - \lambda_B p_B$  when  $\bar{V}(t, n(t)) > 0$ . Applying the infinitesimal generator  $\mathcal{G}$  to  $\bar{V}(t, n(t))$ , and multiplying both sides with  $e^{-\lambda_B t}$ , we get:

$$\begin{aligned} \frac{\partial \bar{V}(t, n(t))}{\partial t} + \lambda_B [\bar{V}(t, n(t) - 1) - \bar{V}(t, n(t))] + \mathcal{G}\Pi(t, n(t)) + \lambda_B p_B &= 0 \\ \frac{\partial \bar{V}(t, n(t))}{\partial t} - \lambda_B \bar{V}(t, n(t)) &= -[\lambda_B \bar{V}(t, n(t) - 1) + \mathcal{G}\Pi(t, n(t)) + \lambda_B p_B] \\ \frac{\partial [\bar{V}(t, n(t))e^{-\lambda_B t}]}{\partial t} &= -[\lambda_B \bar{V}(t, n(t) - 1) + \mathcal{G}\Pi(t, n(t)) + \lambda_B p_B]e^{-\lambda_B t}. \end{aligned}$$

The last differential equation can be solved for  $\bar{V}(t, n(t))$ , provided that  $\bar{V}(t, n(t) - 1)$  is known. Since  $\bar{V}(t, 0) = 0$ , all  $\bar{V}(t, n(t))$  can be solved recursively. The formal procedure is given in the following theorem. We will prove that the  $\bar{V}(t, n(t))$  that is determined by the proposed recursive procedure satisfies conditions (i)-(iv), and is thus equivalent to  $\tilde{V}(t, n(t))$ . Moreover this procedure also determines the latest switching times  $(x_{n(t)})$  for each possible unsold inventory level  $n(t)$ .

**Theorem 5.2.** *For  $1 \leq n(t) \leq M$  and  $\lambda_B > \lambda_i$ , for  $i = 1, 2$ ,  $\bar{V}(t, n(t))$  can be recursively determined by:*

$$\bar{V}(t, n(t)) = \begin{cases} \int_t^T L(s, n(t))e^{-\lambda_B(s-t)} ds & \text{if } t > x_{n(t)} \\ 0 & \text{otherwise,} \end{cases} \quad (22)$$

where

$$\begin{aligned} x_{n(t)} &= \inf \left\{ 0 \leq t \leq T : \int_t^T L(s, n(t))e^{-\lambda_B(s-t)} ds > 0 \right\}, \\ L(t, n(t)) &= \mathcal{G}\Pi(t, n(t)) + \lambda_B p_B + \lambda_B \bar{V}(t, n(t) - 1), \quad 0 \leq t \leq T, \\ \bar{V}(t, 0) &= 0, \quad 0 \leq t \leq T. \end{aligned}$$

What we have shown so far is, for any inventory level  $n = 1, \dots, M$ , there exists a time  $x_n$  such that:  $\bar{V}(t, n) > 0$  if  $t > x_n$  and  $\bar{V}(t, n) = 0$  if  $t \leq x_n$ . Therefore, if the system reaches the  $n$  remaining inventory level at a time  $t \leq x_n$ , then it is optimal to switch immediately. On the other hand if it takes the system longer than  $x_n$  time units to reach the  $n$  remaining inventory level, then it is optimal to delay the switch. Therefore,  $x_n$ 's can

be interpreted as the latest switching time or the switching-time thresholds when  $n$  items are unsold.

Moreover, we showed that these switching-time thresholds  $\{x_n\}$ ,  $n = 1, \dots, M$  are decreasing in unsold inventory  $n$ . Intuitively, as the unsold inventory increases, it is beneficial for the team to delay the switch for more  $t$  values in order to get the advantage of the bundle sales better.

This model captures the essential elements of the problem (i.e. bundling), while allowing sales-to-date to influence the switching decision. The model we use is similar to the one in Feng and Gallego [29], although a crucial feature of ours is that the initial selling period is for the bundled items and after the switch there are two separate processes are active for the individual performances. The limitations of the two-event season and the constant arrival rates will be relaxed in the following section.

## 5.5 Extensions

The dynamic switching problem that is considered so far assumes constant demand rates and a 2-performance selling season. In this section, we relax these two assumptions.

### 5.5.1 $\ell$ -Performances ( $\ell > 2$ ) During the Selling Period

When there are more than two performances on sale, the profit from the individual ticket sales over  $[t, T]$  with  $n(t)$  items available can be expressed as:

$$\Pi(t, n(t)) = \sum_{i=1}^{\ell} p_i E[(N_i(T) - N_i(t)) \wedge n(t)] = \sum_{i=1}^{\ell} \sum_{k=1}^{n(t)} p_i P[N_i(T) - N_i(t) \geq k].$$

It is easy to see that when  $\mathcal{G}$  is applied to  $\Pi(t, n(t))$  we obtain:

$$\mathcal{G}\Pi(t, n(t)) + \lambda_B p_B = (r_B - \sum_{i=1}^{\ell} r_i) + \sum_{i=1}^{\ell} p_i (\lambda_i - \lambda_B) P[N_i(T) - N_i(t) \geq n(t)].$$

As before, we can define a function  $\bar{V}$  which is equivalent to  $\tilde{V}$  as the one defined in Theorem 5.1, the same procedure in Theorem 5.2 is used to calculate it recursively.

**Corollary 5.1.** *When  $\lambda_i \leq \lambda_B$  for each  $i = 1, \dots, \ell$ , and  $r_B > \sum_{i=1}^{\ell} r_i$ , the switching times  $\{x_i\}$   $i = 1, \dots, M$  are decreasing in the remaining inventory.*

### 5.5.2 Time-Dependent Demand Rates

So far we have assumed that the Poisson processes associated with the pre-determined prices have constant demand rates. But in general, the demand rates may change with the remained time, such as they can decrease with time since demand may relate to seat quality. Incorporation of time-dependent demand rates into the formulation will also enable us to indirectly model other aspects such as substitution, where a decrease in the demand rate can be considered to be caused by jockeying customers, or word-of-mouth effect, where demand rates increase with time due to the increase in information about the performance.

We keep the problem setting same as in Section 5.4 and assume that for the Poisson processes:  $N_B(s) 0 \leq s \leq t$  has intensity  $\lambda_B(t)$ ,  $N_1(s) 0 \leq s \leq t$  has intensity  $\lambda_1(t)$ , and  $N_2(s) 0 \leq s \leq t$  has intensity  $\lambda_2(t)$ . Defining  $r_B(t) = \lambda_B(t)p_B$  and  $r_i(t) = \lambda_i(t)p_i$  as the expected revenue rate from the bundled and individual ticket sales of the two performances, we can start our analysis.

To measure the infinitesimal effect of the delay in switching, let us define the infinitesimal generator  $\bar{\mathcal{G}}$  with respect to the Poisson process  $(t, N_B(t))$  for a uniformly bounded function  $g(t, n)$  as:

$$\begin{aligned}\bar{\mathcal{G}}g(t, n) &= \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} E[g(t + \Delta t, n - N_B(\Delta t)) - g(t, n)] \\ &= \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} \sum_{k=0}^{\infty} [g(t + \Delta t, (n - k)^+) - g(t, n)] \frac{(\int_t^{t+\Delta t} \lambda_B(s) ds)^k}{k!} e^{-\int_t^{t+\Delta t} \lambda_B(s) ds} \\ &= \frac{\partial g(t, n)}{\partial t} + \lambda_B(t)[g(t, n - 1) - g(t, n)].\end{aligned}$$

Note that  $\bar{\mathcal{G}}$  is similar to  $\mathcal{G}$  but incorporates the dependence of the demand rate on  $t$ . Applying  $\bar{\mathcal{G}}$  to  $\Pi(t, n(t))$ , we get:

$$\begin{aligned}\bar{\mathcal{G}}\Pi(t, n(t)) + \lambda_B(t)p_B &= [r_B(t) - r_1(t) - r_2(t)] + p_1(\lambda_1(t) - \lambda_B(t))P[N_1(T) - N_1(t) \geq n(t)] \\ &\quad + p_2(\lambda_2(t) - \lambda_B(t))P[N_2(T) - N_2(t) \geq n(t)],\end{aligned}$$

since  $\frac{\partial \sum_{k=1}^{n(t)} P[N_i(T) - N_i(t) \geq k]}{\partial t} = -\lambda_i(t)P[N_i(T) - N_i(t) \leq n(t) - 1]$ . Now we are ready to state the sufficient conditions for the function  $\bar{V}$  which can be calculated recursively, and is identical to  $\tilde{V}$ .

**Theorem 5.3.** *Suppose there exists a function  $\bar{V}(t, n(t))$  such that  $\bar{V}(t, n(t))$  is continuous and differentiable with right continuous derivative in  $[0, T]$  for each fixed  $n(t)$ . In addition, if  $\bar{V}(t, n(t))$  satisfies:*

$$(i) \bar{V}(t, n(t)) \geq 0, \forall 0 \leq t \leq T \text{ and } 0 \leq n(t) \leq M;$$

$$(ii) \bar{V}(T, n(t)) = 0 \text{ for } 0 \leq n(t) \leq M \text{ and } \bar{V}(t, 0) = 0 \text{ for } 0 \leq t \leq T;$$

$$(iii) \bar{V}(t, n(t)) = 0 \Rightarrow \bar{\mathcal{G}}(\bar{V} + \Pi)(t, n(t)) + r_B(t), 0 \leq t \leq T \text{ and } 0 \leq n(t) \leq M;$$

$$(iv) \bar{V}(t, n(t)) > 0 \Rightarrow \bar{\mathcal{G}}(\bar{V} + \Pi)(t, n(t)) + r_B(t), 0 \leq t \leq T \text{ and } 0 \leq n(t) \leq M;$$

then  $\bar{V}(t, n(t)) = \tilde{V}(t, n(t))$ .

The proof is similar to the proof of Theorem 5.1. The main difference is that the demand rates are dependent on time rather than constants. The same procedure described in Theorem 5.2 is used to calculate  $\bar{V}(t, n(t))$  recursively.

**Corollary 5.2.** *When for  $0 \leq t \leq T$ ,  $\lambda_i(t) \leq \lambda_B(t)$  for each  $i = 1, \dots, 2$ , and  $r_B(t) > r_1(t) + r_2(t)$ , the switching times  $\{x_i\}$   $i = 1, \dots, M$  are decreasing in unsold items.*

## 5.6 Computations

In this section, we present the computational analysis illustrating the connection between the problem parameters and the optimal switching times. To calculate the optimal switching times, we need to calculate the revenue potential from delaying the switch,  $\bar{V}(t, n(t))$ , for all inventory levels and times. After dividing the selling period into small time intervals with the size of  $\delta$ , we calculate  $\bar{V}(t, n(t))$ 's recursively starting from  $\bar{V}(T - \delta, 1)$ . The details of the approximation using discrete time intervals (see also [27]) are given below:

For any  $1 \leq n(t) \leq M$ , and  $x_{n(t)} < t < T$  with some  $\delta > 0$  such that  $t + \delta \leq T$ , we have:

$$\begin{aligned}
\bar{V}(t, n(t)) &= \int_t^T L(u, n(t)) e^{-\lambda_B(u-t)} du \\
&= \int_{t+\delta}^T L(u, n(t)) e^{-\lambda_B(u-t)} du + \int_t^{t+\delta} L(u, n(t)) e^{-\lambda_B(u-t)} du \\
&= \int_{t+\delta}^T L(u, n(t)) e^{-\lambda_B(u-(t+\delta))} e^{-\lambda_B \delta} du + \int_t^{t+\delta} L(u, n(t)) e^{-\lambda_B(u-t)} du \\
&\cong \bar{V}(t + \delta, n(t)) e^{-\lambda_B \delta} + \int_t^{t+\delta} L(u, n(t)) e^{-\lambda_B(u-t)} du \\
&\cong \bar{V}(t + \delta, n(t)) e^{-\lambda_B \delta} \\
&\quad + \int_t^{t+\delta} [\mathcal{G}\Pi(u, n(t)) + \lambda_B p_B + \lambda_B \bar{V}(u, n(t) - 1)] e^{-\lambda_B(u-t)} du \\
&\cong \bar{V}(t + \delta, n(t)) e^{-\lambda_B \delta} + (1 - e^{-\lambda_B \delta}) p_B + (1 - e^{-\lambda_B \delta}) \bar{V}(t, n(t) - 1) \\
&\quad + (1 - e^{-\lambda_B \delta}) [\Pi(t, n(t) - 1) - \Pi(t, n(t))] + e^{-\lambda_B \delta} [\Pi(t + \delta, n(t)) - \Pi(t, n(t))].
\end{aligned}$$

Therefore  $\bar{V}(t, n(t))$  can be estimated by:

$$\bar{V}(t, n(t)) \cong (\bar{V} + \Pi)(t + \delta, n(t)) e^{-\lambda_B \delta} + (1 - e^{-\lambda_B \delta}) [p_B + (\bar{V} + \Pi)(t, n(t) - 1)] - \Pi(t, n(t)).$$

If the selling horizon  $T$  is divided into  $K$  (large number of) intervals of length  $\delta$ , we obtain

$$\begin{aligned}
\bar{V}(k\delta, n(t)) &\cong (\bar{V} + \Pi)((k + 1)\delta, n(t)) e^{-\lambda_B \delta} \\
&\quad + (1 - e^{-\lambda_B \delta}) [p_B + (\bar{V} + \Pi)(k\delta, n(t) - 1)] - \Pi(k\delta, n(t)).
\end{aligned}$$

Starting from the end of the selling horizon  $T$ , where  $\bar{V}(T, \cdot) = 0$ , the following algorithm guides computations from inventory level  $n = 1$  to  $M$ .

**Algorithm** Let,

$$\begin{aligned}
\Delta L(k\delta, n(t)) &= (\bar{V} + \Pi)((k + 1)\delta, n(t)) e^{-\lambda_B \delta} \\
&\quad + (1 - e^{-\lambda_B \delta}) [p_B + (\bar{V} + \Pi)(k\delta, n(t) - 1)] - \Pi(k\delta, n(t)).
\end{aligned}$$

- **Step 0:** Initialize  $\bar{V}(T, \cdot) = \bar{V}(K\delta, \cdot) = 0$  for all inventory levels. Set  $n(t) = 1$  and  $k = (K - 1)$ .
- **Step 1:** Calculate  $\Delta L(k\delta, n(t))$ .
- **Step 2:** Set  $\bar{V}(k\delta, n(t)) = (\Delta L(k\delta, n(t)))^+$  and  $k = k - 1$ .

- if  $k \neq -1$  and  $\bar{V}(k\delta, n(t)) \geq 0$ , go to Step 1;
- otherwise set  $\bar{V}(j\delta, n(t)) = 0$  for all  $j < k - 1$  and  $n = n + 1$ .

Consider a team with a 150-ticket stadium facing the problem of selling tickets to one high-demand and one low-demand game during a selling season that lasts 2 months. The demand rates for the games are 50 and 40 seats per month and the prices to be charged to these seats are \$200 and \$50 for high and low-demand games, respectively. If the seats are sold as a bundle with one high and one low-demand seat, the demand rate will be 100 seats for the bundle. Table 8 gives the calculated ten optimal switching times for the case when the bundle is sold with a price of \$220.

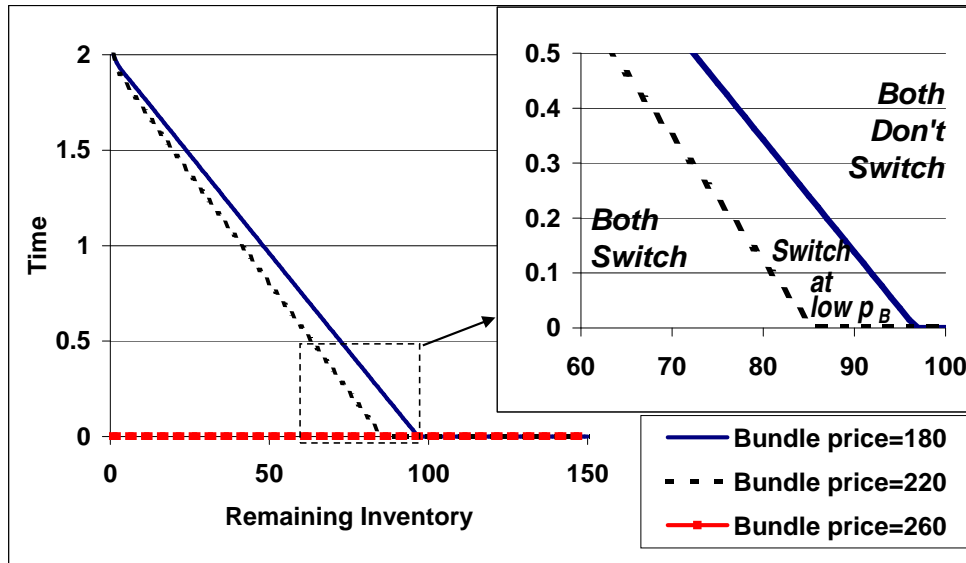
**Table 8:** Selected Optimal Switching Times when  $p_B = 220$

<i>sales</i>	73	72	71	70	69	68	67	66	65	64
remained seats ( $n(t)$ )	77	78	79	80	81	82	83	84	85	86
switch time ( $x_{n(t)}$ )	0.191	0.168	0.145	0.123	0.1	0.078	0.055	0.032	0.01	0

As proved in Theorem 5.2, the optimal switching times are decreasing in unsold inventory  $n$ . Let us consider the case when team has already sold 70 bundles (80 seats remain). In this case the optimal switch time is given as 0.123 months. If the team sold the 70 items more quickly than 0.123 months ( $t < x_{80}$ ), it is optimal to switch before the 71<sup>st</sup> sale arrives since there is no expected revenue potential from delaying the switch further ( $\bar{V}(t, 80) = 0$ ). If the team sold the 70<sup>th</sup> item in bundles after 0.123 months, then they should wait to switch. To illustrate how the switch times are used, it is also beneficial to consider the case when the optimal switch time is 0 with 86 seats leftover to sell in Table 8. Having zero switching times until team sells 64 seats indicates the option of switching should be considered only after the 64<sup>th</sup> sale.

Figure 13 shows the effect of different bundle prices on switching times. As the bundle price increases, the optimal switch threshold decreases for each inventory level, which is intuitive since the team tries to take advantage of high bundle prices by delaying the switch to selling them as singles. If the bundle price is high enough, it may eliminate the switch option altogether such as when  $p_B = 260$ . Figure 13 also illustrates the strategy difference

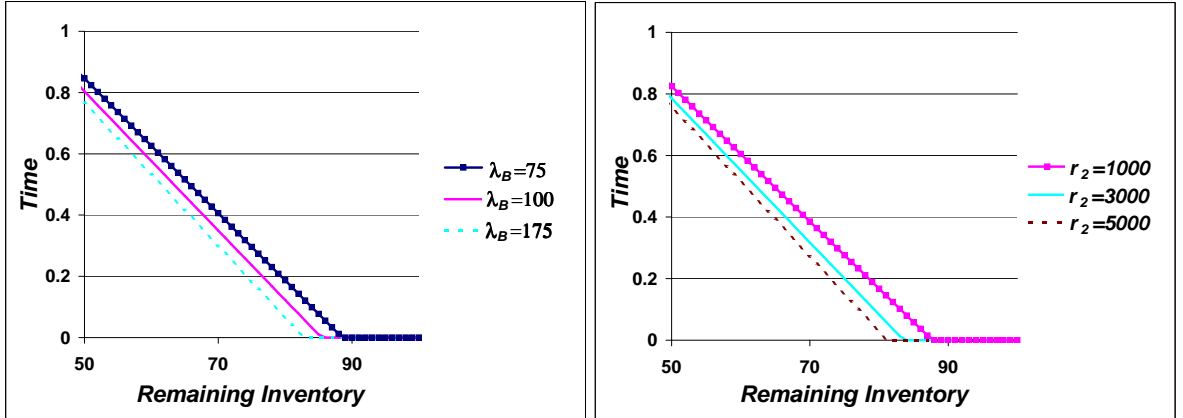
for two different bundle prices ( $p_B = 180$  and  $p_B = 220$ ) with the decision regions.



**Figure 13:** Optimal Switching Times at Different Bundle Prices

Another area of interest is the behavior of the switching thresholds with different demand rates for the bundle. Keeping the prices and demand rates for the single games the same as before, Figure 14(a) illustrates the optimal time thresholds when the bundle is priced at \$220 at different bundle rates. As expected, the time thresholds for switching at each inventory level get smaller (i.e, the time window that requires switching gets smaller) with the increasing demand rates for bundles, which enables the team to take advantage of the high revenue from the bundle for a longer time.

Instead of looking the prices and rates separately, in Figure 14(b) we look at the effect of rate of revenue ( $p \cdot \lambda$ ) for the singles. We consider the case when the rate of revenue from bundles is 22000 with  $p_B = \$220$  and  $\lambda_B = 100$ ; the demand rates for the high and low-demand games are again 50 and 40, respectively. For the singles, we keep the total revenue rate  $r_1 + r_2$  to be constant, so that the relative value of the bundles to the singles does not change. We vary the relative rates of two single events, e.g.,  $r_1=9000$  and  $r_2=3000$  or  $r_1=7000$  and  $r_2=5000$ , keeping the total revenue rate from singles at 12000. Figure



(a) Different Bundle Demand Rates      (b) Different Revenue Rates for Singles

**Figure 14:** Optimal Switching Times for Various Parameters

14(b) shows how the switching time thresholds affected with the change in relative rates of the single events. The switching time thresholds get smaller as the revenue rate from the low-demand events increases (from high-demand events decreases) among the total rate of revenue from the singles. In other words, for a given inventory level, if the low-demand event has higher revenue rate among the singles, team should switch later enabling the team to take advantage of the revenue from the bundle for a longer time.

The final numerical experimentation is performed to see the impact on revenue of deciding the switch time dynamically instead of using a static switch time. The model parameters are the same as the ones that give the optimal switching times in Table 8. For a scenario, we created 100 random sample paths for the arrival of bundled and single ticket customers, and calculate the average revenue when the switch time is decided dynamically and statically for those paths. Figure 15 illustrates that the percentage revenue improvement over the static case by the dynamic switching times changes between 1-2 % when the optimal static switching time (i.e., 1.2) is selected. Another important observation is the reduction in variation when dynamic switch times are used. The revenue values are calculated for the same demand paths, therefore Figure 15 clearly illustrate that the usage of the dynamic switching times improves the value and the predictability of the revenues that will be obtained. Also note that the potential for improvement is higher by the usage of dynamic



switching thresholds when the variation in a scenario is higher.

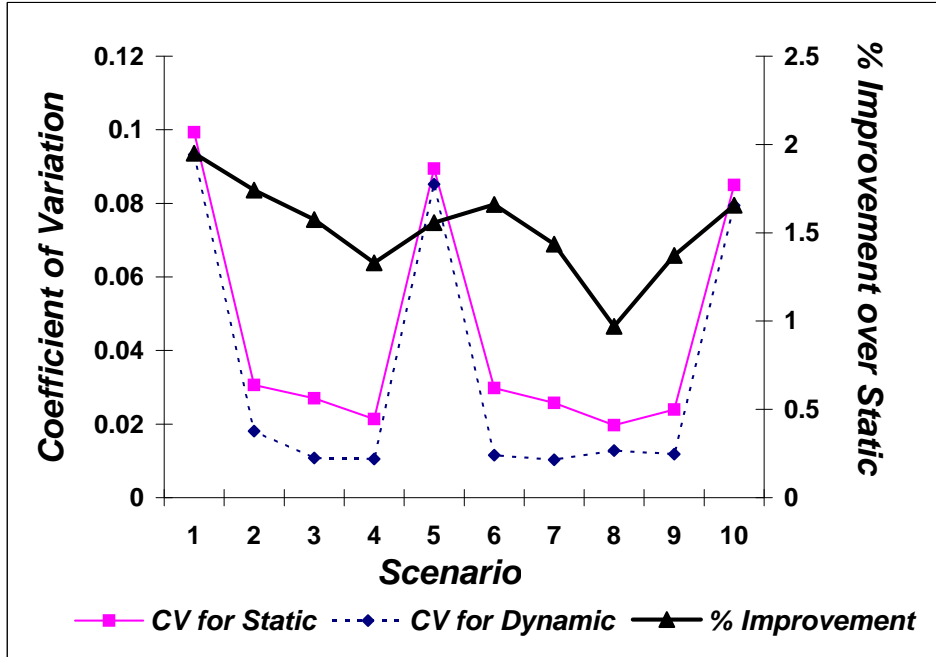


Figure 15: Comparison of Dynamic Decision of Switch Time vs Static

### 5.7 Conclusions

In this chapter, we have studied the problem of switching between selling bundles of tickets to selling individual tickets so as to maximize revenue over a selling season. This is a new problem in revenue management, which is motivated by discussions in the Sports and Entertainment industries, but it may have applications in other industries also. Important characteristics of the problem include that the bundle purchasers are sharing limited capacity with the single-ticket purchaser, and when the switch is made the bundle splits into multiple simultaneous Poisson processes with demand for single tickets. We also focus on the optimal *dynamic* switching time, where the switch time can depend on the state of the system.

Although these characteristics make the problem complicated, we find that the structure yields an optimal policy that is intuitive and easy-to-implement. The optimal time to switch consists of a set of threshold pairs defined by the remaining inventory and the time left in

the horizon. After each sale, the current time is compared to the time threshold for the corresponding remaining inventory to determine if the switch should be made immediately or not. The switching times balance the value of the bundle purchases over the single ticket purchases as well as the probability of future demand arrivals of each. We also generalize these results in several ways including to  $\ell$  events in the horizon or to a demand rate that depends on time, where the same structural results hold. We find that the value of dynamic decision is 1 – 2% over the best static decision, which is comparable to improvements in airline RM.

There are several areas of research to improve the dynamic switching problem. For instance, it would be interesting to explicitly consider how the results change when there are customer diversions between segments. It could also be useful to study how to set the prices of the bundles and single tickets. Models can also be considered that include multiple types of packages (such as full or half-season) in addition to the singles.

There are many useful and interesting research questions that can be analyzed in the Sports and Entertainment industries as a whole. While some of these issues have been studied in other industries, it would be useful to develop decision-making tools specific to the characteristics of S&E, and some of the problems are specific to the S&E industries. For example, some organizations allow consumers to make their own bundles; how should the bundles be formed, how much gain is there from using high-demand events to drive commitment to the bundle, what is the relative value of the increase in demand that may occur with the decrease of commitment to less-popular events? How much would be gained from dynamically adjusting prices in response to demand, and how should prices or inventory of for different categories of seats be determined? These kinds of problems may require OR tools, or OR integrated with economics or marketing. Revenue management in S&E can help to improve the viability of large organizations such as pro-sports, which may face large salary costs to remain competitive, to medium sized organizations (sports or theater), which also struggle to balance costs and revenue, and even small organizations (community theater), where even small improvements may help the organization be sustainable.

## CHAPTER VI

### CONCLUSION

In Chapters 2 and 3, we analyzed the demand management of multi-class customers via stochastic inventory policies by a manufacturer with limited capacity. The customers are segmented into classes based on service level or priority in Chapter 2 and on their acceptance of delay fulfillment in Chapter 3. We showed that modified base-stock policies in the form of  $(S, R, B)$  are optimal for both of the systems in these two chapters, where  $S$  is the optimal order-up-to level,  $R$  denotes the optimal reserve-up-to levels, and  $B$  is the optimal backlog-up-to levels.

These two chapters contribute to the literature on operational models to manage markets with segmented demand, and they show the impact of one kind of flexibility in the production systems. Clearly, the analysis makes assumptions to simplify the problem such as focusing on a single product and using predetermined prices. Yet, these simplifications enabled us to find optimal policies that are easy to understand and implement. Additional research would be helpful to advance the knowledge of how to manage systems with segmented demand. For instance, pricing can be used to control the size of the customer classes in the case of high variability.

In the current business environment, customers can easily share their experiences with each other, informing their decisions about whether or not to do business with a firm or buy that firm's products. The delivery time of an item is not the only factor that affects the customer's decision, but customers may also consider the past performance of a producer. Existing models of leadtime quotation or order acceptance did not capture the impact of past performance on current decisions. In Chapter 4, we consider the optimal demand management of a firm via leadtime quotation when the firm's recent performance of meeting quoted leadtimes impacts the future orders from the customers.

For the infinite capacity case, we find the optimal closed-form expression for leadtime

quotation. We show that the optimal leadtime to quote that accounts for past performance is more conservative (i.e., longer) than the optimal leadtime that ignores it and is always positive unlike the case that ignores service, which means that a firm considering past performance effect would never quote *unethical* leadtimes. When capacity is finite and leadtime is industry-dictated, we determine that the optimal demand acceptance policy does not necessarily have a nice structure, but in some special cases it is convex in the service level of the firm. We develop several heuristics for the general model including a myopic heuristic and one based on balancing the marginal benefit of additional customers, and they perform close to the optimal solution. In our numerical experiments, we find that considering past performance where there is limited capacity can have an impact of more than 2 – 6% in the profit.

Clearly, the usage of the industry-dictated leadtime in the finite capacity case is a limitation in the analysis, but considering that the optimal policy is not well-structured even with this assumption suggests the difficulty of deciding the leadtime to quote. Our research is the first work in the production and leadtime literature to incorporate past performance on customers' decisions; it has shown that service matters significantly in leadtime quotation and demand management, and it is a starting point for many more applications that can be considered in this area.

In the final part of the thesis (Chapter 5), we consider the problem of managing the demand via dynamic timing of the switch between selling bundles of tickets to selling individual tickets. Deciding the optimal time to stop selling season tickets is a new problem in revenue management, which is motivated by discussions in the Sports and Entertainment industries. The main characteristics that differentiate this problem from the ones already studied in revenue management literature include that bundle purchasers are sharing limited capacity with the single-ticket purchaser, and when the switch is made the bundle splits into multiple simultaneous Poisson processes with demand for single tickets. Further, we focus on deciding the switch time *dynamically*.

We show that the optimal time to switch consists of a set of threshold pairs defined by the remaining inventory and the time left in the horizon. After each sale, the current time

is compared to the time threshold for the corresponding remaining inventory to determine if the switch should be made immediately or not. The switching times balance the value of the bundle purchases over the single ticket purchases as well as the probability of future demand arrivals of each. We generalize these results in several ways including considering more than 2 events in the horizon or to a demand rate that depends on time and find the optimal switching times with a similar structure. We also perform numerical experiments to see the effect of deciding the switch time dynamically (instead of using a static switch time decided at the beginning of the selling horizon) on revenue.

There are several areas of research to improve the dynamic switching problem. For instance, it would be interesting to explicitly consider how the results change when there are customer diversions between segments. Offering multiple types of packages (such as full or half-season) in addition to the singles or allowing the sale of bundles and singles simultaneously after the switch are other possible areas of interest to pursue. S&E is a new area of RM, so there may be many others as well.

## APPENDIX A

### PROOFS FOR CHAPTER 2

The following lemma is used in the proof of Theorem 2.1.

**Lemma A.1.** *Given  $g(x, y)$  is jointly concave in  $x$  and  $y$ ,  $G(x) = \max_y g(x, y)$  is a concave function for  $x$ .*

*Proof.* For any  $x_1, x_2 \in R$ , let  $y_1 = \arg \max\{y|g(x_1, y)\}$ ,  $y_2 = \arg \max\{y|g(x_2, y)\}$ . For any  $\lambda \in [0, 1]$ , let  $x_\lambda = \lambda x_1 + (1 - \lambda)x_2$ ,  $y_\lambda = \lambda y_1 + (1 - \lambda)y_2$ . We have  $G(x_\lambda) = \max_y g(x_\lambda, y) \geq g(x_\lambda, y_\lambda) \geq \lambda g(x_1, y_1) + (1 - \lambda)g(x_2, y_2) = \lambda G(x_1) + (1 - \lambda)G(x_2)$ .  $\square$

#### **Proof of Problem Simplification with Nesting for PDS (Lemma 2.1)**

*Proof.* We will show this result by contradiction. Let us start with the first condition,  $(B_t^1 + B_t^2) \cdot R_t^1 = 0$ . Assume that there exists an optimal policy in the form of  $\{(B_t^1 + B_t^2), R_t^1\}$ , where  $(B_t^1 + B_t^2) \cdot R_t^1 > 0$ . We will show that there exists an alternate policy, which is at least as good as and sometimes better than the “optimal” policy, which will contradict the assumption of optimality of the policy where both the reserve inventory decisions and the backlogged order decisions are positive. We consider two main market environments: 1) when the current net revenue from selling out of inventory is better than the future expected profit of an additional unit and 2) when the future expected profit of an additional unit is better than the current net revenue from backlogging.

**Case 1:** Since the current net revenue from selling out of inventory is better than the future marginal expected profit, the alternative policy is saving one less item in the current period.

So the alternate policy is  $\{\overline{B_t^1 + B_t^2}, \overline{R_t^1}\} = \{B_t^1 + B_t^2, R_t^1 - 1\}$ . In both policies, decisions for the second class ( $R_t^2$  and  $B_t^2$ ) are the same. But in the alternate policy, the values of  $S_t^2$  and  $B_t^{2,ef}$  can be higher than the values of the assumed-optimal policy. Let  $V_t$  and  $\overline{V}_t$

be the expected profit starting from period  $t$  under the two policies, respectively. Let us consider the following case:

- When  $D_t^1 - S_t + R_t^1 > B_t^1 + B_t^2 \Rightarrow \overline{B_t^{2,ef}} = B_t^{2,ef} = 0$  and  $\overline{S_t^2} = S_t^2 = 0$

$$\begin{aligned} V_t &= p_t^1(S_t - R_t^1 + B_t^1 + B_t^2) - h_t R_t^1 - \ell_t^1(D_t^1 - S_t + R_t^1 - B_t^1 - B_t^2) - \beta_t^1(B_t^1 + B_t^2) \\ &\quad - \ell_t^2 D_t^2 + J_{t+1}(R_t^1 - B_t^1 - B_t^2) \\ \overline{V}_t &= p_t^1(S_t - R_t^1 + 1 + B_t^1 + B_t^2) - h_t(R_t^1 - 1) - \ell_t^1(D_t^1 - S_t + R_t^1 - 1 - B_t^1 - B_t^2) \\ &\quad - \beta_t^1(B_t^1 + B_t^2) - \ell_t^2 D_t^2 + J_{t+1}(R_t^1 - 1 - B_t^1 - B_t^2) \\ &= V_t + (p_t^1 + \ell_t^1 + h_t) - (J_{t+1}(R_t^1 - B_t^1 - B_t^2) - J_{t+1}(R_t^1 - 1 - B_t^1 - B_t^2)) > V_t \end{aligned}$$

The last inequality follows from the fact that the current net revenue from selling out of inventory is higher than the marginal expected profit from carrying one more unit of inventory forward in this market environment.

**Case 2:** Since the future marginal expected profit is better than the current net revenue from backlogging, promising one less item in the current period is the alternate policy.

So the alternate policy is  $\{\overline{B_t^1 + B_t^2}, \overline{R_t^1}\} = \{B_t^1 + B_t^2 - 1, R_t^1\}$ . In both policies, decisions for the second class ( $R_t^2$  and  $B_t^2$ ) are the same. If we compare  $V_t$  and  $\overline{V}_t$ ,

- When  $D_t^1 - S_t - R_t^1 \geq B_t^1 + B_t^2$ , we will have;

$$\begin{aligned} V_t &= p_t^1(S_t - R_t^1 + B_t^1 + B_t^2) - h_t R_t^1 - \ell_t^1(D_t^1 - S_t + R_t^1 - B_t^1 - B_t^2) - \ell_t^2 D_t^2 \\ &\quad - \beta_t^1(B_t^1 + B_t^2) + J_{t+1}(R_t^1 - B_t^1 - B_t^2) \\ \overline{V}_t &= p_t^1(S_t - R_t^1 + B_t^1 + B_t^2 - 1) - h_t R_t^1 - \ell_t^1(D_t^1 - S_t + R_t^1 - B_t^1 - B_t^2 + 1) - \ell_t^2 D_t^2 \\ &\quad - \beta_t^1(B_t^1 + B_t^2 - 1) + J_{t+1}(R_t^1 - B_t^1 - B_t^2 + 1) \\ &= V_t + (J_{t+1}(R_t^1 - B_t^1 - B_t^2 + 1) - J_{t+1}(R_t^1 - B_t^1 - B_t^2)) - (p_t^1 + \ell_t^1 - \beta_t^1) > V_t \end{aligned}$$

The last inequality follows from the fact that the marginal future expected profit from one more unit of inventory is higher than the current net revenue from backlogging in this market environment.

For the second condition again assume that there exists an optimal policy in the form of  $\{B_t^2, (R_t^1 + R_t^2)\}$  where  $B_t^2 \cdot (R_t^1 + R_t^2) > 0$ . We will show that there exists an alternate policy,

which is at least as good as and sometimes better than the original policy. We consider the same two main market environments as for the first condition.

**Case 1:** Since the current net revenue from selling out of inventory is better than the future marginal expected profit, the alternative policy is saving one less item in the current period.

So the alternate policy is  $\{\overline{B_t^2}, \overline{R_t^1 + R_t^2}\} = \{B_t^2, R_t^1 + R_t^2 - 1\}$ . In both policies, decisions for the first class are the same, namely, the items saved from first-class customers are  $R_t^1$ , and the maximum amount of orders to backlog is  $B_t^1 + B_t^2$ . Let us consider the following case:

- When  $S_t - R_t^1 - R_t^2 \geq D_t^1$ ,  $D_t^2 > S_t^2 - R_t^2 + B_t^2 \Rightarrow \overline{B_t^{2,ef}} = B_t^{2,ef} = B_t^2$  and  $\overline{S_t^2} = S_t^2$

$$V_t = p_t^1 D_t^1 + p_t^2 (S_t^2 - R_t^2 + B_t^2) - h_t (R_t^1 + R_t^2) - \ell_t^2 (D_t^2 - S_t^2 + R_t^2 - B_t^2) - \beta_t^2 B_t^2 \\ + J_{t+1} (R_t^1 + R_t^2 - B_t^2)$$

$$\overline{V}_t = p_t^1 D_t^1 + p_t^2 (S_t^2 - R_t^2 + 1 + B_t^2) - h_t (R_t^1 + R_t^2 - 1) - \ell_t^2 (D_t^2 - S_t^2 + R_t^2 - 1 - B_t^2) \\ - \beta_t^2 B_t^2 + J_{t+1} (R_t^1 + R_t^2 - 1 - B_t^2) \\ = V_t + (p_t^2 + h_t + \ell_t^2) - (J_{t+1} (R_t^1 + R_t^2 - B_t^2) - J_{t+1} (R_t^1 + R_t^2 - 1 - B_t^2)) > V_t$$

The last inequality follows from the fact that the current net revenue from selling out of inventory is higher than the marginal future expected profit from one more unit of inventory in this market environment.

**Case 2:** Since the future marginal expected profit is better than the current net revenue from backlogging, promising one less item in the current period is the alternate policy.

So the alternate policy is  $\{\overline{B_t^2}, \overline{R_t^1 + R_t^2}\} = \{B_t^2 - 1, R_t^1 + R_t^2\}$ . Again, in both policies, decisions for the first class are the same, namely, the items saved from first-class customers are  $R_t^1$ , and the maximum amount of orders to backlog is  $B_t^1 + B_t^2$ . Let us consider the following case:

- When  $S_t - R_t^1 \geq D_t^1 > S_t - R_t^1 - R_t^2$ ,  $D_t^2 \geq B_t^2 \Rightarrow B_t^{2,ef} = B_t^2$ ,  $\overline{B_t^{2,ef}} = B_t^2 - 1$  and



$$\overline{S_t^2} = S_t^2$$

$$\begin{aligned} V_t &= p_t^1 D_t^1 + p_t^2 B_t^2 - h_t(R_t^1 + S_t^2) - \ell_t^2(D_t^2 - B_t^2) - \beta_t^2 B_t^2 + J_{t+1}(S_t^2 + R_t^1 - B_t^2) \\ \overline{V}_t &= p_t^1 D_t^1 + p_t^2(B_t^2 - 1) - h_t(R_t^1 + S_t^2) - \ell_t^2(D_t^2 - B_t^2 + 1) - \beta_t^2(B_t^2 - 1) \\ &\quad + J_{t+1}(S_t^2 + R_t^1 - B_t^2 + 1) \\ &= V_t + (J_{t+1}(R_t^1 + S_t^2 - B_t^2 + 1) - J_{t+1}(R_t^1 + S_t^2 - B_t^2)) - (p_t^2 + \ell_t^2 - \beta_t^2) > V_t \end{aligned}$$

The last inequality follows from the fact that the marginal future expected profit from one more unit of inventory is higher than the current net revenue from backloging in this market environment.

For both of the conditions, the expected profit under the alternative policies is higher than the policy we initially assumed to be optimal in both of the market environments defined at the beginning of the proof, and it can be shown easily that in all other cases in the two market environments, the alternate policies produce exactly the same or higher expected profit as the starting policy. Since alternate policies are at least as good as and sometimes better than the starting policy, a contradiction has been reached.  $\square$

### Proof of Concavity for the Priority Differentiation Strategy (Theorem 2.1)

*Proof.* Let  $j_t(I_t, S_t) = -c_t(S_t - I_t) + G_t(S_t)$ , so  $J_t(I_t) = \max_{S_t: I_t \leq S_t \leq I_t + q_t} j_t(I_t, S_t)$ . We prove by induction.

1. For last period:

$$\begin{aligned} G_T(S_T) &= \int_0^{S_T} p_T^1 \cdot k \cdot d\Phi_T^1(k) + \int_{S_T}^{\infty} p_T^1 \cdot S_T \cdot d\Phi_T^1(k) \\ &\quad + \int_0^{S_T} d\Phi_T^1(k_1) \left( \int_0^{S_T - k_1} k_2 \cdot p_T^2 \cdot d\Phi_T^2(k_2) \right) \\ &\quad + \int_0^{S_T} d\Phi_T^1(k_1) \left( \int_{S_T - k_1}^{\infty} (S_T - k_1) \cdot p_T^2 \cdot d\Phi_T^2(k_2) \right) + \int_0^{S_T} v \cdot (S_T - k) d\Phi_T^{1,2}(k) \end{aligned}$$

Where  $v$  is the salvage value per item at the end of the horizon:  $p_T^1 > p_T^2 > v > 0$

It is clear that the first derivative of  $G_T$  is equal to:

$$G'_T(S_T) = \int_{S_T}^{\infty} p_T^1 d\Phi_T^1(k) + \int_0^{S_T} d\Phi_T^{1,2}(k) \cdot v + \int_0^{S_T} d\Phi_T^1(k_1) \left( \int_{S_T - k_1}^{\infty} p_T^2 d\Phi_T^2(k_2) \right)$$

Now we can check whether  $G'_T(S_T)$  is non-increasing or not:

$$G''_T(S_T) = \phi_T^1(S_T)(p_T^2 - p_T^1) + \phi_T^{1,2}(S_T) \cdot (v - p_T^2)$$

Since  $p_T^1 > p_T^2 > v > 0$ , it is easily seen that  $G''_T(S_T) \leq 0$ , therefore  $G_T(S_T)$  is concave.

2. Given  $t + 1 \leq T$ , assume that  $G_{t+1}(S_T)$  is concave in  $S_T$ , then we can prove that  $j_{t+1}(I_{t+1}, S_{t+1})$  is jointly concave in  $I_{t+1}$  and  $S_{t+1}$  by the following.

For any  $(I_1, S_1), (I_2, S_2) \in \mathfrak{R}^2$ , let  $I_\lambda = \lambda I_1 + (1 - \lambda)I_2$ ,  $S_\lambda = \lambda S_1 + (1 - \lambda)S_2$ . Then,

$$\begin{aligned} j_{t+1}(I_\lambda, S_\lambda) &= -c_{t+1}(S_\lambda - I_\lambda) + G_{t+1}(S_\lambda) \\ &= -c_{t+1}(\lambda S_1 + (1 - \lambda)S_2 - \lambda I_1 - (1 - \lambda)I_2) + G_{t+1}(\lambda S_1 + (1 - \lambda)S_2) \\ &\geq -\lambda c_{t+1}(S_1 - I_1) - (1 - \lambda)c_{t+1}(S_2 - I_2) + \lambda G_{t+1}(S_1) + (1 - \lambda)G_{t+1}(S_2) \\ &= \lambda j_{t+1}(I_1, S_1) + (1 - \lambda)j_{t+1}(I_2, S_2). \end{aligned}$$

So by Lemma A.1,  $J_{t+1}(I_t)$  is concave in  $I_t$ , and as a result  $J'_{t+1}(I_t)$  is non-increasing in  $I_t$ .

3. Next let us prove that  $g_t^1(S_t, R_t^1, R_t^2)$  is quasi-concave in  $R_t^1$  and  $R_t^2$ .

$$\begin{aligned} \frac{\partial g_t^1(S_t, R_t^1, R_t^2)}{\partial R_t^2} &= \begin{cases} 0 & \text{if } S_t \leq R_t^1 + R_t^2 \\ \int_0^{S_t - R_t^1 - R_t^2} d\Phi_t^1(k_1) \left( \int_{S_t - R_t^1 - R_t^2 - k_1}^\infty (J'_{t+1}(R_t^1 + R_t^2) - (p_t^2 + \ell_t^2 + h_t)) d\Phi_t^2(D_t^2) \right) & \text{o.w} \end{cases} \\ \frac{\partial g_t^1(S_t, R_t^1, R_t^2)}{\partial R_t^1} &= \begin{cases} 0 & \text{if } S_t \leq R_t^1 \\ \frac{\partial g_t^1(S_t, R_t^1, R_t^2)}{\partial R_t^2} + \int_{S_t - R_t^1}^\infty [J'_{t+1}(R_t^1) - (p_t^1 + \ell_t^1 + h_t)] d\Phi_t^1(D_t^1) & \text{o.w} \end{cases} \end{aligned}$$

Let us define  $R_t^{1*}$  and  $R_t^{2*}$  as:

$$\begin{aligned} R_t^{1*} &= \max\{I : p_t^1 + \ell_t^1 + h_t \leq J'_{t+1}(I)\} \quad \text{if } p_t^1 + \ell_t^1 + h_t < J'_{t+1}(0) \quad (= 0 \text{ o.w.}) \\ R_t^{1*} + R_t^{2*} &= \max\{I : p_t^2 + \ell_t^2 + h_t \leq J'_{t+1}(I)\} \quad \text{if } p_t^2 + \ell_t^2 + h_t < J'_{t+1}(0) \quad (= 0 \text{ o.w.}). \end{aligned}$$

Thus, we have  $\nabla g_t^1(S_t, R_t^1, R_t^2) \geq [0, 0]^T$  when  $0 \leq R_t^1 \leq R_t^{1*}$  and  $0 \leq R_t^2 \leq R_t^{2*}$ , and  $\nabla g_t^1(S_t, R_t^1, R_t^2) \leq [0, 0]^T$  when  $R_t^1 > R_t^{1*}$  and  $R_t^2 > R_t^{2*}$ ; thus,  $g_t^1(S_t, R_t^1, R_t^2)$  is quasi-concave with respect to  $R_t^1$  and  $R_t^2$ .  $(R_t^{1*}, R_t^{2*})$  is the unique unconstrained optimizer of  $g_t^1(S_t, R_t^1, R_t^2)$ , and it is independent of inventory level  $S_t$ .  $(R_t^{1,c}, R_t^{2,c}) = (\min(R_t^{1*}, (S_t)^+), \min(R_t^{2*}, (S_t)^+))$  maximizes  $g_t^1(S_t, R_t^1, R_t^2)$ , for  $0 \leq R_t^1 \leq (S_t)^+$  and  $0 \leq R_t^2 \leq (S_t)^+$ .

4. Next let us prove that  $g_t^2(S_t, B_t^1, B_t^2)$  is quasi-concave in  $B_t^1$  and  $B_t^2$ .

$$\frac{\partial g_t^2(S_t, B_t^1, B_t^2)}{\partial B_t^1} = \int_{S_t + B_t^1 + B_t^2}^\infty [(p_t^1 + \ell_t^1 - \beta_t^1) - J'_{t+1}(-B_t^1 - B_t^2)] d\Phi_t^1(D_t^1)$$

$$\frac{\partial g_t^2(S_t, B_t^1, B_t^2)}{\partial B_t^2} = \frac{\partial g_t^2(S_t, B_t^1, B_t^2)}{\partial B_t^1} + \int_0^{S_t + B_t^2} d\Phi_t^1(k_1) \left( \int_{S_t + B_t^2 - k_1}^{\infty} [(p_t^2 + \ell_t^2 - \beta_t^2) - J'_{t+1}(-B_t^2)] d\Phi_t^2(D_t^2) \right)$$

Let us define  $B_t^{1*}$  and  $B_t^{2*}$  as:

$$\begin{aligned} B_t^{1*} + B_t^{2*} &= \min\{I : J'_{t+1}(-I) \geq p_t^1 + \ell_t^1 - \beta_t^1\} \quad \text{if } p_t^1 + \ell_t^1 - \beta_t^1 > J'_{t+1}(0) \quad (= 0 \text{ o.w.}) \\ B_t^{2*} &= \min\{I : J'_{t+1}(-I) \geq p_t^2 + \ell_t^2 - \beta_t^2\} \quad \text{if } p_t^2 + \ell_t^2 - \beta_t^2 > J'_{t+1}(0) \quad (= 0 \text{ o.w.}). \end{aligned}$$

Thus, we have  $\nabla g_t^2(S_t, B_t^1, B_t^2) \geq [0, 0]^T$  when  $0 \leq B_t^1 \leq B_t^{1*}$  and  $0 \leq B_t^2 \leq B_t^{2*}$ , and  $\nabla g_t^2(S_t, B_t^1, B_t^2) \leq [0, 0]^T$  when  $B_t^1 > B_t^{1*}$  and  $B_t^2 > B_t^{2*}$ ; thus,  $g_t^2(S_t, B_t^1, B_t^2)$  is quasi-concave with respect to  $B_t^1$  and  $B_t^2$ .  $(B_t^{1*}, B_t^{2*})$  is the unique unconstrained optimizer of  $g_t^2(S_t, B_t^1, B_t^2)$ , and it is independent of inventory level  $S_t$ .  $(B_t^{1,c}, B_t^{2,c}) = (\min(B_t^{1*}, q_{t+1}), \min(B_t^{2*}, q_{t+1}))$  maximizes  $g_t^2(S_t, B_t^1, B_t^2)$ , for  $0 \leq B_t^1 \leq q_{t+1}$  and  $0 \leq B_t^2 \leq q_{t+1}$ .

5. Next let us prove that  $g_t^3(S_t, R_t^2, B_t^1)$  is quasi-concave in  $B_t^1$  and  $R_t^2$ .

$$\begin{aligned} \frac{\partial g_t^3(S_t, R_t^2, B_t^1)}{\partial R_t^2} &= \begin{cases} 0 & \text{if } S_t \leq R_t^2 \\ \int_0^{S_t - R_t^2} d\Phi_t^1(k_1) \left( \int_{S_t - R_t^2 - k_1}^{\infty} [J'_{t+1}(R_t^2) - (p_t^2 + \ell_t^2 + h_t)] d\Phi_t^2(D_t^2) \right) & \text{o.w.} \end{cases} \\ \frac{\partial g_t^3(S_t, R_t^2, B_t^1)}{\partial B_t^1} &= \int_{S_t + B_t^1}^{\infty} [(p_t^1 + \ell_t^1 - \beta_t^1) - J'_{t+1}(-B_t^1)] d\Phi_t^1(k) \end{aligned}$$

Let us define  $R_t^{2*}$  and  $B_t^{1*}$  as:

$$\begin{aligned} R_t^{2*} &= \max\{I : p_t^2 + \ell_t^2 + h_t \leq J'_{t+1}(I)\} \quad \text{if } p_t^2 + \ell_t^2 + h_t < J'_{t+1}(0) \quad (= 0 \text{ o.w.}) \\ B_t^{1*} &= \min\{I : J'_{t+1}(-I) \geq p_t^1 + \ell_t^1 - \beta_t^1\} \quad \text{if } p_t^1 + \ell_t^1 - \beta_t^1 > J'_{t+1}(0) \quad (= 0 \text{ o.w.}). \end{aligned}$$

Thus, we have  $\nabla g_t^3(S_t, R_t^2, B_t^1) \geq [0, 0]^T$  when  $0 \leq R_t^2 \leq R_t^{2*}$  and  $0 \leq B_t^1 \leq B_t^{1*}$ , and  $\nabla g_t^3(S_t, R_t^2, B_t^1) \leq [0, 0]^T$  when  $R_t^2 > R_t^{2*}$  and  $B_t^1 > B_t^{1*}$ ; thus,  $g_t^3(S_t, R_t^2, B_t^1)$  is quasi-concave with respect to  $R_t^2$  and  $B_t^1$ .  $(R_t^{2*}, B_t^{1*})$  is the unique unconstrained optimizer of  $g_t^3(S_t, R_t^2, B_t^1)$ , and it is independent of inventory level  $S_t$ .  $(R_t^{2,c}, B_t^{1,c}) = (\min(R_t^{2*}, S_t), \min(B_t^{1*}, q_{t+1}))$  maximizes  $g_t^3(S_t, R_t^2, B_t^1)$ , for  $0 \leq R_t^2 \leq (S_t)^+$  and  $0 \leq B_t^1 \leq q_{t+1}$ .

6. Let us prove the concavity of  $G_t^1(S_t)$  with respect to  $S_t$ , where  $G_t^1(S_t) = g_t^1(S_t, R_t^{1,c}, R_t^{2,c})$ .

We will consider  $G_t^1(S_t)$  in five cases:

Case I:  $S_t \leq R_t^{1*}$

The profit-to-go after production in this case is:  $G_t^1(S_t) = -h_t S_t + J_{t+1}(S_t) -$

$\int_0^\infty \ell_t^1 k d\Phi_t^1(k) - \int_0^\infty \ell_t^2 k d\Phi_t^2(k)$ , so its first derivative is  $G_t^{1'}(S_t) = -h_t + J'_{t+1}(S_t)$ . Thus, it is clear that  $G_t^{1'}(S_t)$  is non-increasing since:  $G_t^{1''}(S_t) = J''_{t+1}(S_t)$  and  $J''_{t+1}(S_t) \leq 0$ .

Case II:  $R_t^{1*} + \epsilon < S_t \leq R_t^{1*} + R_t^{2*}$

The first derivative of profit-to-go after production function is:

$$G_t^{1'}(S_t) = \int_{S_t - R_t^{1*}}^\infty p_t^1 d\Phi_t^1(k) - h_t \int_0^{S_t - R_t^{1*}} d\Phi_t^1(k) + \int_0^{S_t - R_t^{1*}} J'_{t+1}(S_t - k) d\Phi_t^1(k) \\ + \int_{S_t - R_t^{1*}}^\infty \ell_t^1 d\Phi_t^1(k).$$

Thus, it is clear that  $G_t^{1'}(S_t)$  is non-increasing since:

$$G_t^{1''}(S_t) = \int_0^{S_t - R_t^{1*}} J''_{t+1}(S_t - k) d\Phi_t^1(k) + \phi_t^1(S_t - R_t^{1*})(J'_{t+1}(R_t^{1*}) - (p_t^1 + h_t + \ell_t^1));$$

also,  $J''_{t+1}(S_t) \leq 0$ , and  $p_t^1 + h_t + \ell_t^1 = J'_{t+1}(R_t^{1*})$  due to the  $R_t^{1*}$  decisions.

Case III:  $S_t = R_t^{1*} + \epsilon$

It is clear that  $G_t^{1'}(S_t)$  is non-increasing in this case since:

$$G_t^{1''}(R_t^{1*} + \epsilon) = (1 - \phi_t^1(0))(p_t^1 + h_t + \ell_t^1 - J'_{t+1}(R_t^{1*})) \\ + \phi_t^1(0)(J'_{t+1}(R_t^{1*} + \epsilon) - J'_{t+1}(R_t^{1*}));$$

also,  $J'_{t+1}(S_t)$  is non-increasing and  $p_t^1 + h_t + \ell_t^1 = J'_{t+1}(R_t^{1*})$ .

Case IV:  $R_t^{1*} + R_t^{2*} + \epsilon < S_t$

In this case the first derivative of the profit-to-go after production is equal to:

$$G_t^{1'}(S_t) = \int_{S_t - R_t^{1*}}^\infty p_t^1 d\Phi_t^1(k) - \int_{S_t - R_t^{1*} - R_t^{2*}}^{S_t - R_t^{1*}} h_t d\Phi_t^1(k) \\ + \int_{S_t - R_t^{1*} - R_t^{2*}}^{S_t - R_t^{1*}} J'_{t+1}(S_t - k) d\Phi_t^1(k) \\ + \int_0^{S_t - R_t^{1*} - R_t^{2*}} d\Phi_t^1(k_1) \left( \int_{S_t - R_t^{1*} - R_t^{2*} - k_1}^\infty p_t^2 d\Phi_t^2(k_2) \right) \\ - \int_0^{S_t - R_t^{1*} - R_t^{2*}} d\Phi_t^1(k_1) \left( \int_0^{S_t - R_t^{1*} - R_t^{2*} - k_1} h_t d\Phi_t^2(k_2) \right) + \int_{S_t - R_t^{1*}}^\infty \ell_t^1 d\Phi_t^1(k) \\ + \int_0^{S_t - R_t^{1*} - R_t^{2*}} d\Phi_t^1(k_1) \left( \int_0^{S_t - R_t^{1*} - R_t^{2*} - k_1} J'_{t+1}(S_t - k_2 - k_1) d\Phi_t^2(k_2) \right) \\ + \int_0^{S_t - R_t^{1*} - R_t^{2*}} d\Phi_t^1(k_1) \left( \int_{S_t - R_t^{1*} - R_t^{2*} - k_1}^\infty \ell_t^2 d\Phi_t^2(k_2) \right).$$

The second derivative is:

$$G_t^{1''}(S_t) = \phi_t^1(S_t - R_t^{1*} - R_t^{2*})(p_t^2 + h_t + \ell_t^2 - J'_{t+1}(R_t^{1*} + R_t^{2*})) \\ + \int_0^{S_t - R_t^{1*} - R_t^{2*}} d\Phi_t^1(k_1) \phi_t^2(S_t - R_t^{1*} - R_t^{2*} - k_1)(J'_{t+1}(R_t^{1*} + R_t^{2*}) - (p_t^2 + h_t + \ell_t^2)) \\ + \int_0^{S_t - R_t^{1*} - R_t^{2*}} d\Phi_t^1(k_1) \left( \int_0^{S_t - R_t^{1*} - R_t^{2*} - k_1} (J''_{t+1}(S_t - k_2 - k_1) d\Phi_t^2(k_2)) \right) \\ + \int_{S_t - R_t^{1*} - R_t^{2*}}^{S_t - R_t^{1*}} J''_{t+1}(S_t - k) d\Phi_t^1(k) + \phi_t^1(S_t - R_t^{1*})(J'_{t+1}(R_t^{1*}) - (p_t^1 + h_t + \ell_t^1)).$$

Thus, it is clear that  $G_t^{1'}(S_t)$  is non-increasing since  $J''_{t+1}(S_t) \leq 0$ ,  $p_t^1 + h_t + \ell_t^1 = J'_{t+1}(R_t^{1*})$  due to the  $R_t^{1*}$  decision, and  $J'_{t+1}(R_t^{1*} + R_t^{2*}) = p_t^2 + h_t + \ell_t^2$  due to the  $R_t^{1*} + R_t^{2*}$  decision.

Case V:  $R_t^{1*} + R_t^{2*} + \epsilon = S_t$

In this case, the second derivative is:

$$\begin{aligned} G_t^{1''}(S_t) &= \phi_t^1(0)(1 - \phi_t^2(0))(p_t^2 + h_t + \ell_t^2 - J'_{t+1}(R_t^{1*} + R_t^{2*})) \\ &\quad + (J'_{t+1}(R_t^{1*} + R_t^{2*} + \epsilon) - J'_{t+1}(R_t^{1*} + R_t^{2*}))\phi_t^1(0)\phi_t^2(0) \\ &\quad + \phi_t^1(R_t^{2*})(J'_{t+1}(R_t^{1*}) - (p_t^1 + h_t + \ell_t^1)) \\ &\quad + \int_{\epsilon}^{R_t^{2*}} (J'_{t+1}(R_t^{1*} + R_t^{2*} + \epsilon - k) - J'_{t+1}(R_t^{1*} + R_t^{2*} - k))d\Phi_t^1(k). \end{aligned}$$

It is clear that  $G_t^1(S_t)$  is non-increasing in this case since  $J'_{t+1}(S_t)$  is non-increasing,  $p_t^1 + h_t + \ell_t^1 = J'_{t+1}(R_t^{1*})$  due to the  $R_t^{1*}$  decision, and  $J_{t+1}(R_t^{1*} + R_t^{2*}) = p_t^2 + h_t + \ell_t^2$  due to the  $R_t^{1*} + R_t^{2*}$  decision.

7. Let us prove the concavity of  $G_t^2(S_t)$  with respect to  $S_t$ , where  $G_t^2(S_t) = g_t^2(S_t, B_t^{1,c}, B_t^{2,c})$ .

The first derivative of profit-to-go after production is equal to:

$$\begin{aligned} G_t^{2'}(S_t) &= \int_{S_t+B_t^{1*}+B_t^{2*}}^{\infty} p_t^1 d\Phi_t^1(k) + \int_0^{S_t+B_t^{2*}} d\Phi_t^1(k_1) \left( \int_{S_t+B_t^{2*}-k_1}^{\infty} p_t^2 d\Phi_t^2(k_2) \right) \\ &\quad + \int_0^{S_t+B_t^{2*}} J'_{t+1}(S_t - k) d\Phi_t^T(k) + \int_{S_t+B_t^{2*}}^{S_t+B_t^{1*}+B_t^{2*}} J'_{t+1}(S_t - k) d\Phi_t^1(k) \\ &\quad + \int_{S_t+B_t^{1*}+B_t^{2*}}^{\infty} \ell_t^1 d\Phi_t^1(k) + \int_0^{S_t+B_t^{2*}} d\Phi_t^1(k_1) \left( \int_{S_t+B_t^{2*}-k_1}^{\infty} \ell_t^2 d\Phi_t^2(k_2) \right) \\ &\quad + \int_{S_t}^{S_t+B_t^{1*}+B_t^{2*}} \beta_t^1 d\Phi_t^1(k) + \int_0^{S_t} d\Phi_t^1(k_1) \left( \int_{S_t-k_1}^{S_t-k_1+B_t^{2*}} \beta_t^2 d\Phi_t^2(k_2) \right) \\ &\quad - \int_{S_t}^{S_t+B_t^{2*}} d\Phi_t^1(k_1) \left( \int_{S_t+B_t^{2*}-k_1}^{\infty} \beta_t^2 d\Phi_t^2(k_2) \right) - h_t \int_0^{S_t} d\Phi_t^T(k). \end{aligned}$$

The second derivative is:

$$\begin{aligned} G_t^{2''}(S_t) &= \phi_t^1(S_t + B_t^{1*} + B_t^{2*})(J'_{t+1}(-B_t^{1*} - B_t^{2*})) - (p_t^1 + \ell_t^1 - \beta_t^1) \\ &\quad + \int_0^{S_t+B_t^{2*}} J''_{t+1}(S_t - k) d\Phi_t^T(k) + \phi_t^1(S_t)(\beta_t^2 - \beta_t^1) \\ &\quad - (h_t + \beta_t^2)\phi_t^T(S_t) + \phi_t^1(S_t + B_t^{2*})(p_t^2 + \ell_t^2 - \beta_t^2 - J'_{t+1}(-B_t^{2*})) \\ &\quad + \phi_t^T(S_t + B_t^{2*})(J'_{t+1}(-B_t^{2*}) - (p_t^2 + \ell_t^2 - \beta_t^2)) \\ &\quad + \int_{S_t+B_t^{2*}}^{S_t+B_t^{1*}+B_t^{2*}} J''_{t+1}(S_t - k) d\Phi_t^1(k). \end{aligned}$$

Thus, it is clear that  $G_t^2(S_t)$  is non-increasing since  $J''_{t+1}(S_t) \leq 0$ ;  $p_t^2 + \ell_t^2 - \beta_t^2 = J'_{t+1}(-B_t^{2*})$  due to the  $B_t^{2*}$  decision;  $J_{t+1}(-B_t^{1*} - B_t^{2*}) = p_t^1 + \ell_t^1 - \beta_t^1$  due to the  $B_t^{1*} + B_t^{2*}$  decision, and  $\beta_t^1 \geq \beta_t^2$ .

8. Let us prove the concavity of  $G_t^3(S_t)$  with respect to  $S_t$ , where  $G_t^3(S_t) = g_t^3(S_t, R_t^{2,c}, B_t^{1,c})$ .

We will consider  $G_t^3(S_t)$  in three cases:

Case I:  $S_t \leq R_t^{2*}$

So the first and second derivatives of the profit-to-go after production in this case are

equal to:

$$G_t^{3'}(S_t) = \int_{S_t+B_t^{1*}}^{\infty} p_t^1 d\Phi_t^1(k) - \int_0^{S_t} h_t d\Phi_t^1(k) + \int_0^{S_t+B_t^{1*}} J'_{t+1}(S_t-k) d\Phi_t^1(k) \\ + \int_{S_t+B_t^{1*}}^{\infty} \ell_t^1 d\Phi_t^1(k) + \int_{S_t}^{S_t+B_t^{1*}} \beta_t^1 d\Phi_t^1(k).$$

$$G_t^{3''}(S_t) = \phi_t^1(S_t+B_t^{1*})(J'_{t+1}(-B_t^{1*}) - (p_t^1 + \ell_t^1 - \beta_t^1)) - (h_t + \beta_t^1)\phi_t^1(S_t) \\ + \int_0^{S_t+B_t^{1*}} J''_{t+1}(S_t-k) d\Phi_t^1(k).$$

Thus, it is clear that  $G_t^{3'}(S_t)$  is non-increasing since  $J''_{t+1}(S_t) \leq 0$  and  $p_t^1 + \ell_t^1 - \beta_t^1 = J'_{t+1}(-B_t^{1*})$  due to the  $B_t^{1*}$  decisions.

Case II:  $R_t^{2*} + \epsilon < S_t$

Then the first and second derivatives of the profit-to-go after production are equal to:

$$G_t^{3'}(S_t) = \int_{S_t+B_t^{1*}}^{\infty} p_t^1 d\Phi_t^1(k) - h_t \int_{S_t-R_t^{2*}}^{S_t} d\Phi_t^1(k) + \int_{S_t+B_t^{1*}}^{\infty} \ell_t^1 d\Phi_t^1(k) \\ + \int_0^{S_t-R_t^{2*}} d\Phi_t^1(k_1) \left( \int_{S_t-R_t^{2*}-k_1}^{\infty} d\Phi_t^2(k_2) p_t^2 \right) + \int_{S_t}^{S_t+B_t^{1*}} \beta_t^1 d\Phi_t^1(k) \\ - \int_0^{S_t-R_t^{2*}} d\Phi_t^1(k_1) \left( \int_0^{S_t-R_t^{2*}-k_1} d\Phi_t^2(k_2) h_t \right) + \int_{S_t-R_t^{2*}}^{S_t+B_t^{1*}} J'_{t+1}(S_t-k) d\Phi_t^1(k) \\ + \int_0^{S_t-R_t^{2*}} d\Phi_t^1(k_1) \left( \int_0^{S_t-R_t^{2*}-k_1} J'_{t+1}(S_t-k_2-k_1) d\Phi_t^2(k_2) \right) \\ + \int_0^{S_t-R_t^{2*}} d\Phi_t^1(k_1) \left( \int_{S_t-R_t^{2*}-k_1}^{\infty} \ell_t^2 d\Phi_t^2(k_2) \right).$$

$$G_t^{3''}(S_t) = \phi_t^1(S_t-R_t^{2*})(p_t^2 + h_t + \ell_t^2 - J'_{t+1}(R_t^{2*})) \\ + \phi_t^1(S_t+B_t^{1*})(J'_{t+1}(-B_t^{1*}) - (p_t^1 + \ell_t^1 - \beta_t^1)) \\ + \int_0^{S_t-R_t^{2*}} d\Phi_t^1(k_1) \phi_t^2(S_t-R_t^{2*}-k_1)(J'_{t+1}(R_t^{2*}) - (p_t^2 + h_t + \ell_t^2)) \\ + \int_0^{S_t-R_t^{2*}} d\Phi_t^1(k_1) \left( \int_0^{S_t-R_t^{2*}-k_1} J''_{t+1}(S_t-k_2-k_1) d\Phi_t^2(k_2) \right) \\ + \int_{S_t-R_t^{2*}}^{S_t+B_t^{1*}} J''_{t+1}(S_t-k) d\Phi_t^1(k) - (h_t + \beta_t^1)\phi_t^1(S_t).$$

Thus, it is clear that  $G_t^{3'}(S_t)$  is non-increasing since:  $J''_{t+1}(S_t) \leq 0$ ;  $p_t^1 + \ell_t^1 - \beta_t^1 = J'_{t+1}(-B_t^{1*})$  due to the  $B_t^{1*}$  decision, and  $p_t^2 + h_t + \ell_t^2 = J'_{t+1}(R_t^{2*})$  due to the  $R_t^{2*}$  decision.

Case III:  $R_t^{2*} + \epsilon = S_t$  and  $B_t^{1*} \leq q_{t+1}$

The second derivative of  $G_t^3(S_t)$  is:

$$G_t^{3''}(S_t) = \phi_t^1(0)(1 - \phi_t^2(0))(p_t^2 + h_t + \ell_t^2 - J'_{t+1}(R_t^{2*})) - \beta_t^1 \phi_t^1(R_t^{2*}) \\ + \phi_t^1(R_t^{2*} + B_t^{1*})(J'_{t+1}(-B_t^{1*}) - (p_t^1 + \ell_t^1 - \beta_t^1)) \\ + (J'_{t+1}(R_t^{2*} + \epsilon) - J'_{t+1}(R_t^{2*}))\phi_t^1(0)\phi_t^2(0) - h_t \phi_t^1(R_t^{2*} + \epsilon) \\ + \int_{\epsilon}^{R_t^{2*} + B_t^{1*}} (J'_{t+1}(R_t^{2*} + \epsilon - k) - J'_{t+1}(R_t^{2*} - k)) d\Phi_t^1(k).$$

We know that  $G_t^{3'}(S_t)$  is non-increasing in this case since  $J'_{t+1}(S_t)$  is non-increasing;

$p_t^1 + \ell_t^1 - \beta_t^1 = J'_{t+1}(-B_t^{1*})$  due to the  $B_t^{1*}$  decision, and  $p_t^2 + h_t + \ell_t^2 = J'_{t+1}(R_t^{2*})$  due to the  $R_t^{2*}$  decision.

9. Let us prove the concavity of  $G_t(S_t)$ .

In each period, we must be in one of the following cases, which are independent of the  $S_t$  values:

- If  $J'_{t+1}(0) > p_t^1 + \ell_t^1 + h_t$ , we have  $R_t^{1*} \geq 0$ ,  $R_t^{2*} \geq 0$ ,  $B_t^{1*} = 0$ , and  $B_t^{2*} = 0$ , therefore  $R_t^{1,c} \geq 0$ ,  $R_t^{2,c} \geq 0$ ,  $B_t^{1,c} = 0$ , and  $B_t^{2,c} = 0$ ; thus, we have  $G_t(S_t) = G_t^1(S_t)$ .
- If  $p_t^2 + \ell_t^2 - \beta_t^2 > J'_{t+1}(0)$ , we have  $B_t^{1*} \geq 0$ ,  $B_t^{2*} \geq 0$ ,  $R_t^{1*} = 0$ , and  $R_t^{2*} = 0$ , therefore  $B_t^{1,c} \geq 0$ ,  $B_t^{2,c} \geq 0$ ,  $R_t^{1,c} = 0$ , and  $R_t^{2,c} = 0$ ; thus, we have  $G_t(S_t) = G_t^2(S_t)$ .
- If  $J'_{t+1}(0) > p_t^2 + \ell_t^2 + h_t$ , and  $J'_{t+1}(0) < p_t^1 - \beta_t^1 + \ell_t^1$ , we have  $R_t^{2*} \geq 0$ ,  $R_t^{1*} = 0$ ,  $B_t^{1*} \geq 0$ , and  $B_t^{2*} = 0$ , therefore  $R_t^{1,c} = 0$ ,  $R_t^{2,c} \geq 0$ ,  $B_t^{1,c} \geq 0$ , and  $B_t^{2,c} = 0$ ; thus, we have  $G_t(S_t) = G_t^3(S_t)$ .

So in each period,  $G_t(S_t)$  reduces to some function that is concave. Therefore  $G_t(S_t)$  is concave.

□

## APPENDIX B

### PROOFS FOR CHAPTER 3

#### Proof of Lemma 3.1

*Proof.* We will show this result by contradiction. Assume that there exists an optimal policy in the form of  $\{B_t, (R_t^1 + R_t^2)\}$  where  $B_t \cdot (R_t^1 + R_t^2) > 0$ . We will show that there exists an alternate policy that is at least as good as and sometimes better than the assumed optimal policy, which will contradict the optimality of the assumed policy where both the reserve inventory decisions and the backlogging availability decision are positive. We consider two main market environments: 1) when the current net revenue from selling out of inventory is better than the future expected profit of an additional unit and 2) when the future profit of an additional unit is better than the current net revenue from backlogging.

**Case 1:** Since the current net revenue from selling out of inventory is better than the future expected profit of an additional unit, the alternative policy is saving one item less in the current period.

So the alternate policy is;  $\{\overline{B}_t, \overline{R}_t^1 + \overline{R}_t^2\} = \{B_t, R_t^1 + R_t^2 - 1\}$ . In both policies, The decision for the first class is the same, namely, the items saved from first class customers is  $R_t^1$ . Let  $V_t$  and  $\overline{V}_t$  be the expected profit starting from period  $t$  under the two policies, respectively. Let us consider two cases:

- Case 1.1:  $S_t^2 \geq R_t^2, S_t^2 - R_t^2 < D_t^2 \leq S_t^2 - R_t^2 + B_t \Rightarrow S_t^2 > \overline{R}_t^2, S_t^2 - \overline{R}_t^2 \leq D_t^2 < S_t^2 - \overline{R}_t^2 + \overline{B}_t$

$$V_t = p_t^1 D_t^1 + p_t^2 D_t^2 - h_t R_t^2 - h_t R_t^1 - \beta_t^2 (D_t^2 - S_t^2 + R_t^2) + J_{t+1}^{TDS}(R_t^1 + S_t^2 - D_t^2)$$

$$\overline{V}_t = p_t^1 D_t^1 + p_t^2 D_t^2 - h_t (R_t^2 - 1) - h_t R_t^1 - \beta_t^2 (D_t^2 - S_t^2 + R_t^2 - 1)$$

$$+ J_{t+1}^{TDS}(R_t^1 + S_t^2 - D_t^2) = V_t + h_t + \beta_t^2 > V_t$$



- Case 1.2:  $S_t^2 \geq R_t^2$ ,  $D_t^2 > S_t^2 - R_t^2 + B_t \Rightarrow S_t^2 > \overline{R}_t^2$ ,  $D_t^2 \geq S_t^2 - \overline{R}_t^2 + \overline{B}_t$

$$V_t = p_t^1 D_t^1 + p_t^2 (S_t^2 - R_t^2 + B_t) - h_t R_t^2 - h_t^1 R_t^1 - \ell_t^2 (D_t^2 - S_t^2 + R_t^2 - B_t) - \beta_t^2 B_t \\ + J_{t+1}^{TDS}(R_t^1 + R_t^2 - B_t)$$

$$\overline{V}_t = p_t^1 D_t^1 + p_t^2 (S_t^2 - R_t^2 + 1 + B_t) - h_t (R_t^2 - 1) - h_t^1 R_t^1 - \ell_t^2 (D_t^2 - S_t^2 + R_t^2 - 1 - B_t) \\ - \beta_t^2 B_t + J_{t+1}^{TDS}(R_t^1 + R_t^2 - B_t - 1) \\ = V_t + (p_t^2 + h_t + \ell_t^2) - [J_{t+1}^{TDS}(R_t^1 + R_t^2 - B_t) - J_{t+1}^{TDS}(R_t^1 + R_t^2 - B_t - 1)] > V_t$$

The last inequality follows from the fact that the current net revenue from selling out of inventory is better than the future expected profit of an additional unit in this market setting.

**Case 2:** Since the future profit of an additional unit is better than the current net revenue from backlogging, promising one item less in the current period is the alternate policy.

So the alternate policy is;  $\{\overline{B}_t, \overline{R}_t^1 + \overline{R}_t^2\} = \{B_t - 1, R_t^1 + R_t^2\}$ . Again, in both policies,  $R_t^1$  decision for the first class is same. Let us compare  $V_t$  and  $\overline{V}_t$  under the following three cases:

- Case 2.1:  $D_t^1 \geq S_t - R_t^1$  and  $D_t^2 \geq B_t \Rightarrow D_t^2 > \overline{B}_t$

$$V_t = p_t^1 (S_t - R_t^1) + p_t^2 B_t - h_t R_t^1 - \ell_t^1 (D_t^1 - S_t + R_t^1) - \ell_t^2 (D_t^2 - B_t) - \beta_t^2 B_t \\ + J_{t+1}^{TDS}(R_t^1 - B_t)$$

$$\overline{V}_t = p_t^1 (S_t - R_t^1) + p_t^2 (B_t - 1) - h_t R_t^1 - \ell_t^1 (D_t^1 - S_t + R_t^1) - \ell_t^2 (D_t^2 - B_t + 1) \\ - \beta_t^2 (B_t - 1) + J_{t+1}^{TDS}(R_t^1 - B_t + 1) \\ = V_t + [J_{t+1}^{TDS}(R_t^1 - B_t + 1) - J_{t+1}^{TDS}(R_t^1 - B_t)] - (p_t^2 + \ell_t^2 - \beta_t^2) > V_t$$

- Case 2.2:  $S_t - R_t^1 > D_t^1 > S_t - R_t^1 - R_t^2$  and  $D_t^2 \geq B_t \Rightarrow D_t^2 > \overline{B}_t$

$$V_t = p_t^1 D_t^1 + p_t^2 B_t - h_t S_t^2 - h_t R_t^1 - \ell_t^2 (D_t^2 - B_t) - \beta_t^2 B_t + J_{t+1}^{TDS}(S_t^2 + R_t^1 - B_t)$$

$$\overline{V}_t = p_t^1 D_t^1 + p_t^2 (B_t - 1) - h_t S_t^2 - h_t R_t^1 - \ell_t^2 (D_t^2 - B_t + 1) - \beta_t^2 (B_t - 1) \\ + J_{t+1}^{TDS}(S_t^2 + R_t^1 - B_t + 1)$$

$$= V_t + [J_{t+1}^{TDS}(S_t^2 + R_t^1 - B_t + 1) - J_{t+1}^{TDS}(S_t^2 + R_t^1 - B_t)] - (p_t^2 + \ell_t^2 - \beta_t^2) > V_t$$

- Case 2.3:  $S_t^2 \geq R_t^2$ ,  $D_t^2 \geq S_t^2 - R_t^2 + B_t \Rightarrow S_t^2 \geq \overline{R}_t^2$ ,  $D_t^2 > S_t^2 - \overline{R}_t^2 + \overline{B}_t$

$$V_t = p_t^1 D_t^1 + p_t^2 (S_t^2 - R_t^2 + B_t) - h_t R_t^2 - h_t R_t^1 - \ell_t^2 (D_t^2 - S_t^2 + R_t^2 - B_t) - \beta_t^2 B_t \\ + J_{t+1}^{TDS}(R_t^1 + R_t^2 - B_t)$$

$$\overline{V}_t = p_t^1 D_t^1 + p_t^2 (S_t^2 - R_t^2 + B_t - 1) - h_t R_t^2 - h_t R_t^1 - \ell_t^2 (D_t^2 - S_t^2 + R_t^2 - B_t + 1) \\ - \beta_t^2 (B_t - 1) + J_{t+1}^{TDS}(R_t^1 + R_t^2 - B_t + 1)$$

$$= V_t + [J_{t+1}^{TDS}(R_t^1 + R_t^2 - B_t + 1) - J_{t+1}^{TDS}(R_t^1 + R_t^2 - B_t)] - (p_t^2 + \ell_t^2 - \beta_t^2) > V_t$$

The last inequalities in all three of the cases follows from the fact that the future profit of an additional unit is better than the current net revenue from backloging in this market setting.

The expected profit under the alternative policies is higher than the policy we initially assumed to be optimal in both of the market environments defined at the beginning of the proof, and it can be shown easily that in all other cases in the two market environments, the alternate policies produce exactly the same expected profit. Since alternate policies are at least as good as and sometimes better than the starting policy, a contradiction has been reached.

□

### Proof of Lemma 3.2

*Proof.* By contradiction, assume that there is an optimal policy with  $R_t \cdot B_t > 0$  for some period  $t$ . Let  $\overline{R}_t = R_t - 1$  and  $\overline{B}_t = B_t - 1$  be the alternative policy, and let  $V_t$  and  $\overline{V}_t$  be the expected profit starting from period  $t$  under the two policies respectively. We compare the two policies in the following three cases:

- Case 1:  $D_t^{1,2} \leq S_t - R_t$ , hence  $D_t^{1,2} < S_t - \overline{R}_t$ .

$$V_t = p_t^2 D_t^{1,2} - h_t (S_t - D_t^{1,2}) + J_{t+1}^{CSS}(S_t - D_t^{1,2}) = \overline{V}_t$$

- Case 2:  $S_t - R_t + B_t/\alpha_t > D_t^{1,2} > S_t - R_t$ , hence  $S_t - \overline{R}_t + \overline{B}_t/\alpha_t \geq D_t^{1,2} \geq S_t - \overline{R}_t$ .

$$\begin{aligned}
V_t &= p_t^2(S_t - R_t + \lfloor \alpha_t(D_t^{1,2} - S_t + R_t) \rfloor) - h_t R_t - \beta_t^2 \lfloor \alpha_t(D_t^{1,2} - S_t + R_t) \rfloor \\
&\quad - \ell_t^1 \lceil (1 - \alpha_t)(D_t^{1,2} - S_t + R_t) \rceil + J_{t+1}^{CSS}(R_t - \lfloor \alpha_t(D_t^{1,2} - S_t + R_t) \rfloor) \\
\overline{V}_t &= p_t^2(S_t - R_t + 1 + \lfloor \alpha_t(D_t^{1,2} - S_t + R_t - 1) \rfloor) - h_t(R_t - 1) \\
&\quad - \beta_t^2 \lfloor \alpha_t(D_t^{1,2} - S_t + R_t - 1) \rfloor - \ell_t^1 \lceil (1 - \alpha_t)(D_t^{1,2} - S_t + R_t - 1) \rceil \\
&\quad + J_{t+1}^{CSS}(R_t - 1 - \lfloor \alpha_t(D_t^{1,2} - S_t + R_t - 1) \rfloor)
\end{aligned}$$

If  $\lfloor \alpha_t(D_t^{1,2} - S_t + R_t - 1) \rfloor = \lfloor \alpha_t(D_t^{1,2} - S_t + R_t) \rfloor$ , then  $\lceil (1 - \alpha_t)(D_t^{1,2} - S_t + R_t - 1) \rceil = \lceil (1 - \alpha_t)(D_t^{1,2} - S_t + R_t) \rceil - 1$ . We have,

$$\overline{V}_t = V_t + p_t^2 + h_t + \ell_t^1 - J_{t+1}^{CSS}(R_t - \lfloor \alpha_t(D_t^{1,2} - S_t + R_t) \rfloor) + J_{t+1}^{CSS}(R_t - 1 - \lfloor \alpha_t(D_t^{1,2} - S_t + R_t) \rfloor)$$

Since  $D_t^{1,2} < S_t - R_t + B_t/\alpha_t$ , a new demand from class 2 will be accepted, which means  $p_t^2 + h_t + \ell_t^2 + J_{t+1}^{CSS}(R_t - 1 - \lfloor \alpha_t(D_t^{1,2} - S_t + R_t) \rfloor) \geq J_{t+1}^{CSS}(R_t - \lfloor \alpha_t(D_t^{1,2} - S_t + R_t) \rfloor)$ . Thus  $\overline{V}_t \geq V_t$ .

Otherwise,  $\lfloor \alpha_t(D_t^{1,2} - S_t + R_t - 1) \rfloor = \lfloor \alpha_t(D_t^{1,2} - S_t + R_t) \rfloor - 1$ , then  $\lceil (1 - \alpha_t)(D_t^{1,2} - S_t + R_t - 1) \rceil = \lceil (1 - \alpha_t)(D_t^{1,2} - S_t + R_t) \rceil$ . We have,  $\overline{V}_t = V_t + h_t + \beta_t^2 \geq V_t$ .

- Case 3:  $D_t^{1,2} \geq S_t - R_t + B_t/\alpha_t$ , hence  $D_t^{1,2} > S_t - \overline{R}_t + \overline{B}_t/\alpha_t$ .

$$\begin{aligned}
V_t &= p_t^2(S_t - R_t + B_t) - h_t R_t - \beta_t^2 B_t - \ell_t^1(1 - \alpha_t)B_t/\alpha_t - \ell_t(D_t^{1,2} - S_t + R_t - B_t/\alpha_t) \\
&\quad + J_{t+1}^{CSS}(R_t - B_t) \\
\overline{V}_t &= V_t + h_t + \beta_t^2 + (1 - \alpha_t)(\ell_t^1 - \ell_t^2) \geq V_t
\end{aligned}$$

The expected profit under the alternative policy is always greater or equal to that under the current policy, which incurs a contradiction.  $\square$

### Proof of Concavity Results (Theorem 3.1)

**Lemma B.1.** *Given  $g(x, y)$  is jointly concave in  $x$  and  $y$ ,  $G(x) = \max_y g(x, y)$  is a concave function for  $x$ .*

*Proof.* For any  $x_1, x_2 \in R$ , let  $y_1 = \arg \max\{y|g(x_1, y)\}$ ,  $y_2 = \arg \max\{y|g(x_2, y)\}$ . For any  $\lambda \in [0, 1]$ , let  $x_\lambda = \lambda x_1 + (1 - \lambda)x_2$ ,  $y_\lambda = \lambda y_1 + (1 - \lambda)y_2$ . We have  $G(x_\lambda) = \max_y g(x_\lambda, y) \geq g(x_\lambda, y_\lambda) \geq \lambda g(x_1, y_1) + (1 - \lambda)g(x_2, y_2) = \lambda G(x_1) + (1 - \lambda)G(x_2)$ .  $\square$

### For the Time Differentiation Strategy

In the proof below, the *TDS* superscript is omitted from the expected profit functions to increase readability.

*Proof.* Let  $j_t(I_t, S_t) = -c_t(S_t - I_t) + G_t(S_t)$ , so  $J_t(I_t) = \max_{S_t: I_t \leq S_t \leq I_t + q_t} j_t(I_t, S_t)$ . We prove by induction.

1. For period  $t = T$ , we have  $B_T = 0$ ,  $R_T^1 = R_T^2 = 0$ , and  $J_{T+1}(I_T) = v \cdot I_T$ .

$$\begin{aligned} G_T(S_T) &= \int p_T^1 \min(D_T^1, S_T) d\Phi_T^1(D_T^1) \\ &\quad + \iint p_T^2 \min(D_T^2, [S_T - D_T^1]^+) d\Phi_T^1(D_T^1) d\Phi_T^2(D_T^2) \\ &\quad + \iint v \cdot \max(0, S_T - D_T^1 - D_T^2) d\Phi_T^1(D_T^1) d\Phi_T^2(D_T^2) \\ &= \int_0^{S_T} p_T^1 \cdot k \cdot d\Phi_T^1(k) + \int_{S_T}^\infty p_T^1 \cdot S_T \cdot d\Phi_T^1(k) \\ &\quad + \int_0^{S_T} d\Phi_T^1(k_1) \left( \int_0^{S_T - k_1} k_2 \cdot p_T^2 \cdot d\Phi_T^2(k_2) \right) \\ &\quad + \int_0^{S_T} d\Phi_T^1(k_1) \left( \int_{S_T - k_1}^\infty (S_T - k_1) \cdot p_T^2 \cdot d\Phi_T^2(k_2) \right) + \int_0^{S_T} v \cdot (S_T - k) d\Phi_T^{1,2}(k) \end{aligned}$$

Where  $v$  is the salvage value per item at the end of the horizon:  $p_T^1 > p_T^2 > v > 0$

It is clear that its first derivative is equal to:

$$G'_T(S_T) = \int_{S_T}^\infty p_T^1 d\Phi_T^1(k) + \int_0^{S_T} d\Phi_T^{1,2}(k) \cdot v + \int_0^{S_T} d\Phi_T^1(k_1) \left( \int_{S_T - k_1}^\infty p_T^2 d\Phi_T^2(k_2) \right)$$

Know we can check whether  $G'_T(S_T)$  is non-increasing or not:

$$G''_T(S_T) = \phi_T^1(S_T)(p_T^2 - p_T^1) + \phi_T^{1,2}(S_T) \cdot (v - p_T^2)$$

Since  $p_T^1 > p_T^2 > v > 0$ , it is easily seen that  $G''_T(S_T) \leq 0$ , therefore  $G_T(S_T)$  is concave.

2. Given  $t + 1 \leq T$ , assume that  $G_{t+1}(S_t)$  is concave in  $S_t$ , then we can prove that  $j_{t+1}(I_{t+1}, Y_{t+1})$  is jointly concave in  $I_{t+1}$  and  $Y_{t+1}$  by the following.

For any  $(I_1, Y_1), (I_2, Y_2) \in \mathfrak{R}^2$ , let  $I_\lambda = \lambda I_1 + (1 - \lambda)I_2$ ,  $Y_\lambda = \lambda Y_1 + (1 - \lambda)Y_2$ . Then,

$$\begin{aligned} j_{t+1}(I_\lambda, Y_\lambda) &= -c_{t+1}(Y_\lambda - I_\lambda) + G_{t+1}(Y_\lambda) \\ &= -c_{t+1}(\lambda Y_1 + (1 - \lambda)Y_2 - \lambda I_1 - (1 - \lambda)I_2) + G_{t+1}(\lambda Y_1 + (1 - \lambda)Y_2) \\ &\geq -\lambda c_{t+1}(Y_1 - I_1) - (1 - \lambda)c_{t+1}(Y_2 - I_2) + \lambda G_{t+1}(Y_1) + (1 - \lambda)G_{t+1}(Y_2) \\ &= \lambda j_{t+1}(I_1, Y_1) + (1 - \lambda)j_{t+1}(I_2, Y_2). \end{aligned}$$

So by Lemma B.1,  $J_{t+1}(I_t)$  is concave in  $I_t$ , and as a result  $J'_{t+1}(I_t)$  is non-increasing in  $I_t$ .

3. Next let us prove that  $g_t^{TDS}(S_t, R_t^1, R_t^2, 0)$  is quasi-concave in  $R_t^1$  and  $R_t^2$ .

$$\frac{\partial g_t^{TDS}(S_t, R_t^1, R_t^2, 0)}{\partial R_t^2} = \begin{cases} 0 & \text{if } S_t \leq R_t^1 + R_t^2 \\ \int_0^{S_t - R_t^1 - R_t^2} d\Phi_t^1(k_1) \left( \int_{S_t - R_t^1 - R_t^2 - k_1}^{\infty} (J'_{t+1}(R_t^1 + R_t^2) - (p_t^2 + \ell_t^2 + h_t)) d\Phi_t^2(D_t^2) \right) & \text{o.w.} \end{cases}$$

$$\frac{\partial g_t^{TDS}(S_t, R_t^1, R_t^2, 0)}{\partial R_t^1} = \begin{cases} 0 & \text{if } S_t \leq R_t^1 \\ \frac{\partial g_t^1(S_t, R_t^1, R_t^2)}{\partial R_t^1} + \int_{S_t - R_t^1}^{\infty} [J'_{t+1}(R_t^1) - (p_t^1 + \ell_t^1 + h_t)] d\Phi_t^1(D_t^1) & \text{o.w.} \end{cases}$$

Let us define  $R_t^{1*}$  and  $R_t^{2*}$  as:

$$R_t^{1*} = \begin{cases} \max\{I : p_t^1 + \ell_t^1 + h_t \leq J'_{t+1}(I)\} & \text{if } p_t^1 + \ell_t^1 + h_t < J'_{t+1}(0) \\ 0 & \text{otherwise.} \end{cases}$$

$$R_t^{1*} + R_t^{2*} = \begin{cases} \max\{I : p_t^2 + \ell_t^2 + h_t \leq J'_{t+1}(I)\} & \text{if } p_t^2 + \ell_t^2 + h_t < J'_{t+1}(0) \\ 0 & \text{otherwise.} \end{cases}$$

Thus we have  $\nabla g_t^{TDS}(S_t, R_t^1, R_t^2, 0) \geq [0, 0]^T$  when  $0 \leq R_t^1 \leq R_t^{1*}$  and  $0 \leq R_t^2 \leq R_t^{2*}$ , and

$\nabla g_t^{TDS}(S_t, R_t^1, R_t^2, 0) \leq [0, 0]^T$  when  $R_t^1 > R_t^{1*}$  and  $R_t^2 > R_t^{2*}$ ; thus,  $g_t^{TDS}(S_t, R_t^1, R_t^2, 0)$  is quasi-concave with respect to  $R_t^1$  and  $R_t^2$ .  $(R_t^{1*}, R_t^{2*})$  is the unique unconstrained optimizer of  $g_t^{TDS}(S_t, R_t^1, R_t^2, 0)$  and it is independent of inventory level  $S_t$ .  $(R_t^{1,c}, R_t^{2,c}) = (\min(R_t^{1*}, (S_t)^+), \min(R_t^{2*}, (S_t)^+))$  maximizes  $g_t^{TDS}(S_t, R_t^1, R_t^2, 0)$ , for  $0 \leq R_t^1 \leq (S_t)^+$  and  $0 \leq R_t^2 \leq (S_t)^+$ .

4. Next let us prove that  $g_t^{TDS}(S_t, 0, 0, B_t)$  is quasi-concave in  $B_t$ .

$$\frac{\partial g_t^{TDS}(S_t, 0, 0, B_t)}{\partial B_t} = \int_0^{S_t} d\Phi_t^1(k_1) \left( \int_{S_t + B_t - k_1}^{\infty} [p_t^2 + \ell_t^2 - \beta_t^2 - J'_{t+1}(-B_t)] d\Phi_t^2(D_t^2) \right) + \int_{S_t}^{\infty} d\Phi_t^1(k_1) \left( \int_{B_t}^{\infty} [p_t^2 + \ell_t^2 - \beta_t^2 - J'_{t+1}(-B_t)] d\Phi_t^2(D_t^2) \right)$$

Let us define  $B_t^*$  as:

$$B_t^* = \begin{cases} \min\{I : J'_{t+1}(-I) \geq p_t^2 + \ell_t^2 - \beta_t^2\} & \text{if } p_t^2 + \ell_t^2 - \beta_t^2 > J'_{t+1}(0) \\ 0 & \text{otherwise.} \end{cases}$$

Thus we have  $g_t'^{TDS}(S_t, 0, 0, B_t) \geq 0$  when  $0 \leq B_t \leq B_t^*$ , and  $g_t'^{TDS}(S_t, 0, 0, B_t) \leq 0$  when  $B_t > B_t^*$ ; thus,  $g_t^{TDS}(S_t, 0, 0, B_t)$  is quasi-concave with respect to  $B_t$ .  $B_t^*$  is the unique unconstrained optimizer of  $g_t^{TDS}(S_t, 0, 0, B_t)$  and it is independent of inventory level  $S_t$ .  $B_t^c = \min(B_t^*, q_{t+1})$  maximizes  $g_t^{TDS}(S_t, 0, 0, B_t)$ , for  $0 \leq B_t \leq q_{t+1}$ .

5. Let us prove the concavity of  $G_t^R(S_t)$  with respect to  $S_t$ , where

$G_t^R(S_t) = g_t^{TDS}(S_t, R_t^{1,c}, R_t^{2,c}, 0)$ . We will consider  $G_t^R(S_t)$  in five cases:

Case I:  $S_t \leq R_t^{1*}$

The profit-to-go after production in this case is:  $G_t^R(S_t) = -h_t S_t + J_{t+1}(S_t) - \int_0^\infty \ell_t^1 k d\Phi_t^1(k) - \int_0^\infty \ell_t^2 k d\Phi_t^2(k)$ , so its first derivative is  $G_t^{R'}(S_t) = -h_t + J_{t+1}'(S_t)$ . Thus it is clear that  $G_t^{R'}(S_t)$  is non-increasing since:  $G_t^{R''}(S_t) = J_{t+1}''(S_t)$  and  $J_{t+1}''(S_t) \leq 0$ .

Case II:  $R_t^{1*} + \epsilon < S_t \leq R_t^{1*} + R_t^{2*}$

The profit-to-go after production in this case is:

$$\begin{aligned} G_t^R(S_t) &= \int_0^{S_t - R_t^{1*}} k p_t^1 d\Phi_t^1(k) + \int_{S_t - R_t^{1*}}^\infty p_t^1 (S_t - R_t^{1*}) d\Phi_t^1(k) - h_t R_t^{1*} \\ &\quad - \int_0^{S_t - R_t^{1*}} (S_t - R_t^{1*} - k) h_t d\Phi_t^1(k) + \int_0^{S_t - R_t^{1*}} J_{t+1}(S_t - k) d\Phi_t^1(k) \\ &\quad + \int_{S_t - R_t^{1*}}^\infty J_{t+1}(R_t^{1*}) d\Phi_t^1(k) - \int_0^\infty \ell_t^2 k d\Phi_t^2(k) \\ &\quad - \int_{S_t - R_t^{1*}}^\infty \ell_t^1 (k - (S_t - R_t^{1*})) d\Phi_t^1(k). \end{aligned}$$

Its first derivative is:

$$\begin{aligned} G_t^{R'}(S_t) &= \int_{S_t - R_t^{1*}}^\infty p_t^1 d\Phi_t^1(k) - h_t \int_0^{S_t - R_t^{1*}} d\Phi_t^1(k) + \int_0^{S_t - R_t^{1*}} J_{t+1}'(S_t - k) d\Phi_t^1(k) \\ &\quad + \int_{S_t - R_t^{1*}}^\infty \ell_t^1 d\Phi_t^1(k). \end{aligned}$$

Thus it is clear that  $G_t^{R'}(S_t)$  is non-increasing since:

$$G_t^{R''}(S_t) = \int_0^{S_t - R_t^{1*}} J_{t+1}''(S_t - k) d\Phi_t^1(k) + \phi_t^1(S_t - R_t^{1*}) (J_{t+1}'(R_t^{1*}) - (p_t^1 + h_t + \ell_t^1));$$

also,  $J_{t+1}''(S_t) \leq 0$ , and  $p_t^1 + h_t + \ell_t^1 = J_{t+1}'(R_t^{1*})$  due to the  $R_t^{1*}$  decisions.

Case III:  $S_t = R_t^{1*} + \epsilon$

Now it is clear that  $G_t^{R'}(S_t)$  is non-increasing in this case since:

$$\begin{aligned} G_t^{R''}(R_t^{1*} + \epsilon) &= \int_\epsilon^\infty p_t^1 d\Phi_t^1(k) - \int_0^\epsilon h_t d\Phi_t^1(k) + \int_0^\epsilon J_{t+1}'(R_t^{1*} + \epsilon - k) d\Phi_t^1(k) \\ &\quad + \int_\epsilon^\infty \ell_t^1 d\Phi_t^1(k) - J_{t+1}'(R_t^{1*}) + h_t \\ &= (1 - \phi_t^1(0))(p_t^1 + h_t + \ell_t^1 - J_{t+1}'(R_t^{1*})) \\ &\quad + \phi_t^1(0)(J_{t+1}'(R_t^{1*} + \epsilon) - J_{t+1}'(R_t^{1*})); \end{aligned}$$

also,  $J_{t+1}'(S_t)$  is non-increasing and  $p_t^1 + h_t + \ell_t^1 = J_{t+1}'(R_t^{1*})$ .

Case IV:  $R_t^{1*} + R_t^{2*} + \epsilon < S_t$

So the profit-to-go after production in this case is:

$$\begin{aligned}
G_t^R(S_t) &= \int_0^{S_t - R_t^{1*}} k \cdot p_t^1 d\Phi_t^1(k) + \int_{S_t - R_t^{1*}}^\infty p_t^1 (S_t - R_t^{1*}) d\Phi_t^1(k) - h_t R_t^{1*} \\
&\quad - \int_0^{S_t - R_t^{1*} - R_t^{2*}} R_t^{2*} h_t d\Phi_t^1(k) - \int_{S_t - R_t^{1*} - R_t^{2*}}^{S_t - R_t^{1*}} h_t (S_t - R_t^{1*} - k) d\Phi_t^1(k) \\
&\quad + \int_0^{S_t - R_t^{1*} - R_t^{2*}} d\Phi_t^1(k_1) \left( \int_0^{S_t - R_t^{1*} - R_t^{2*} - k_1} k_2 \cdot p_t^2 d\Phi_t^2(k_2) \right) \\
&\quad + \int_0^{S_t - R_t^{1*} - R_t^{2*}} d\Phi_t^1(k_1) \left( \int_{S_t - R_t^{1*} - R_t^{2*} - k_1}^\infty p_t^2 (S_t - R_t^{1*} - R_t^{2*} - k_1) d\Phi_t^2(k_2) \right) \\
&\quad - h_t \int_0^{S_t - R_t^{1*} - R_t^{2*}} d\Phi_t^1(k_1) \left( \int_0^{S_t - R_t^{1*} - R_t^{2*} - k_1} (S_t - R_t^{1*} - R_t^{2*} - k_1 - k_2) d\Phi_t^2(k_2) \right) \\
&\quad + \int_0^{S_t - R_t^{1*} - R_t^{2*}} d\Phi_t^1(k_1) \left( \int_0^{S_t - R_t^{1*} - R_t^{2*} - k_1} J_{t+1}(S_t - k_1 - k_2) d\Phi_t^2(k_2) \right) \\
&\quad + \int_0^{S_t - R_t^{1*} - R_t^{2*}} d\Phi_t^1(k_1) \left( \int_{S_t - R_t^{1*} - R_t^{2*} - k_1}^\infty J_{t+1}(R_t^{1*} + R_t^{2*}) d\Phi_t^2(k_2) \right) \\
&\quad + \int_{S_t - R_t^{1*} - R_t^{2*}}^{S_t - R_t^{1*}} J_{t+1}(S_t - k) d\Phi_t^1(k) - \int_{S_t - R_t^{1*} - R_t^{2*}}^\infty d\Phi_t^1(k_1) \left( \int_0^\infty k_2 \ell_t^2 d\Phi_t^2(k_2) \right) \\
&\quad - \int_{S_t - R_t^{1*}}^\infty \ell_t^1 (k - (S_t - R_t^{1*})) d\Phi_t^1(k) + \int_{S_t - R_t^{1*}}^\infty J_{t+1}(R_t^{1*}) d\Phi_t^1(k) \\
&\quad - \int_0^{S_t - R_t^{1*} - R_t^{2*}} d\Phi_t^1(k_1) \left( \int_{S_t - R_t^{1*} - R_t^{2*} - k_1}^\infty (k_2 - S_t + R_t^{1*} + R_t^{2*} + k_1) \ell_t^2 d\Phi_t^2(k_2) \right).
\end{aligned}$$

Its first derivative is equal to:

$$\begin{aligned}
G_t^{R'}(S_t) &= \int_{S_t - R_t^{1*}}^\infty p_t^1 d\Phi_t^1(k) - \int_{S_t - R_t^{1*} - R_t^{2*}}^{S_t - R_t^{1*}} h_t d\Phi_t^1(k) \\
&\quad + \int_{S_t - R_t^{1*} - R_t^{2*}}^{S_t - R_t^{1*}} J'_{t+1}(S_t - k) d\Phi_t^1(k) \\
&\quad + \int_0^{S_t - R_t^{1*} - R_t^{2*}} d\Phi_t^1(k_1) \left( \int_{S_t - R_t^{1*} - R_t^{2*} - k_1}^\infty p_t^2 d\Phi_t^2(k_2) \right) \\
&\quad - \int_0^{S_t - R_t^{1*} - R_t^{2*}} d\Phi_t^1(k_1) \left( \int_0^{S_t - R_t^{1*} - R_t^{2*} - k_1} h_t d\Phi_t^2(k_2) \right) + \int_{S_t - R_t^{1*}}^\infty \ell_t^1 d\Phi_t^1(k) \\
&\quad + \int_0^{S_t - R_t^{1*} - R_t^{2*}} d\Phi_t^1(k_1) \left( \int_0^{S_t - R_t^{1*} - R_t^{2*} - k_1} J'_{t+1}(S_t - k_2 - k_1) d\Phi_t^2(k_2) \right) \\
&\quad + \int_0^{S_t - R_t^{1*} - R_t^{2*}} d\Phi_t^1(k_1) \left( \int_{S_t - R_t^{1*} - R_t^{2*} - k_1}^\infty \ell_t^2 d\Phi_t^2(k_2) \right).
\end{aligned}$$

The second derivative is:

$$\begin{aligned}
G_t^{R''}(S_t) &= \phi_t^1(S_t - R_t^{1*} - R_t^{2*}) (p_t^2 + h_t + \ell_t^2 - J'_{t+1}(R_t^{1*} + R_t^{2*})) \\
&\quad + \int_0^{S_t - R_t^{1*} - R_t^{2*}} d\Phi_t^1(k_1) \phi_t^2(S_t - R_t^{1*} - R_t^{2*} - k_1) (J'_{t+1}(R_t^{1*} + R_t^{2*}) - (p_t^2 + h_t + \ell_t^2)) \\
&\quad + \int_0^{S_t - R_t^{1*} - R_t^{2*}} d\Phi_t^1(k_1) \left( \int_0^{S_t - R_t^{1*} - R_t^{2*} - k_1} (J''_{t+1}(S_t - k_2 - k_1) d\Phi_t^2(k_2)) \right) \\
&\quad + \int_{S_t - R_t^{1*} - R_t^{2*}}^{S_t - R_t^{1*}} J''_{t+1}(S_t - k) d\Phi_t^1(k) + \phi_t^1(S_t - R_t^{1*}) (J'_{t+1}(R_t^{1*}) - (p_t^1 + h_t + \ell_t^1)).
\end{aligned}$$

Thus it is clear that  $G_t^{R'}(S_t)$  is non-increasing since  $J''_{t+1}(S_t) \leq 0$ ,  $p_t^1 + h_t + \ell_t^1 = J'_{t+1}(S_t)$  due to the  $R_t^{1*}$  decision, and  $J'_{t+1}(R_t^{1*} + R_t^{2*}) = (p_t^2 + h_t + \ell_t^2)$  due to the  $R_t^{2*}$  decision.

Case V:  $R_t^{1*} + R_t^{2*} + \epsilon = S_t$

In this case, the second derivative is:

$$\begin{aligned}
G_t^{R''}(S_t) &= \int_{R_t^{2*} + \epsilon}^{\infty} p_t^1 d\Phi_t^1(k) - h_t \int_{\epsilon}^{R_t^{2*} + \epsilon} d\Phi_t^1(k) + h_t \int_0^{R_t^{2*}} d\Phi_t^1(k) \\
&\quad + \int_{\epsilon}^{R_t^{2*} + \epsilon} J'_{t+1}(R_t^{1*} + R_t^{2*} + \epsilon - k) d\Phi_t^1(k) \\
&\quad + \int_0^{\epsilon} d\Phi_t^1(k_1) \left( \int_0^{\epsilon - k_1} J'_{t+1}(R_t^{1*} + R_t^{2*} + \epsilon - k_1 - k_2) d\Phi_t^2(k_2) \right) \\
&\quad - \int_0^{R_t^{2*}} J'_{t+1}(R_t^{1*} + R_t^{2*} - k) d\Phi_t^1(k) - \int_0^{\epsilon} d\Phi_t^1(k_1) \left( \int_0^{\epsilon - k_1} h_t d\Phi_t^2(k_2) \right) \\
&\quad - \int_{R_t^{2*}}^{\infty} p_t^1 d\Phi_t^1(k) + \int_{R_t^{2*} + \epsilon}^{\infty} \ell_t^1 d\Phi_t^1(k) - \int_{R_t^{2*}}^{\infty} \ell_t^1 d\Phi_t^1(k) \\
&\quad + \int_0^{\epsilon} d\Phi_t^1(k_1) \left( \int_{\epsilon - k_1}^{\infty} \ell_t^2 d\Phi_t^2(k_2) \right) + \int_0^{\epsilon} d\Phi_t^1(k_1) \left( \int_{\epsilon - k_1}^{\infty} p_t^2 d\Phi_t^2(k_2) \right) \\
&= \phi_t^1(0)(1 - \phi_t^2(0))(p_t^2 + h_t + \ell_t^2 - J'_{t+1}(R_t^{1*} + R_t^{2*})) \\
&\quad + (J'_{t+1}(R_t^{1*} + R_t^{2*} + \epsilon) - J'_{t+1}(R_t^{1*} + R_t^{2*}))\phi_t^1(0)\phi_t^2(0) \\
&\quad + \phi_t^1(R_t^{2*})(J'_{t+1}(R_t^{1*}) - (p_t^1 + h_t + \ell_t^1)) \\
&\quad + \int_{\epsilon}^{R_t^{2*}} (J'_{t+1}(R_t^{1*} + R_t^{2*} + \epsilon - k) - J'_{t+1}(R_t^{1*} + R_t^{2*} - k)) d\Phi_t^1(k).
\end{aligned}$$

It is clear that  $G_t^{R'}(S_t)$  is non-increasing in this case since  $J'_{t+1}(S_t)$  is non-increasing,  $p_t^1 + h_t + \ell_t^1 = J'_{t+1}(R_t^{1*})$  due to the  $R_t^{1*}$  decision, and  $J_{t+1}(R_t^{1*} + R_t^{2*}) = p_t^2 + h_t + \ell_t^2$  due to the  $R_t^{2*}$  decision.

6. Let us prove the concavity of  $G_t^B(S_t)$  with respect to  $S_t$ , where

$G_t^B(S_t) = g_t^{TDS}(S_t, 0, 0, B_t^c)$ . We will consider  $G_t^B(S_t)$  in two cases:

Case I:  $B_t^* \leq q_{t+1}$

Its first derivative is equal to:

$$\begin{aligned}
G_t^{B'}(S_t) &= \int_{S_t}^{\infty} (p_t^1 + \ell_t^1) d\Phi_t^1(k) + \int_0^{S_t} d\Phi_t^1(k_1) \left( \int_{S_t + B_t^* - k_1}^{\infty} p_t^2 d\Phi_t^2(k_2) \right) \\
&\quad + \int_0^{S_t} d\Phi_t^1(k) \left( \int_{S_t - k_1}^{S_t - k_1 + B_t^*} J'_{t+1}(S_t - k_1 - k_2) d\Phi_t^2(k_2) \right) \\
&\quad - \int_0^{S_t} h_t d\Phi_t^{1,2}(k) + \int_0^{S_t} d\Phi_t^1(k_1) \left( \int_{S_t + B_t^* - k_1}^{\infty} \ell_t^2 d\Phi_t^2(k_2) \right) \\
&\quad + \int_0^{S_t} d\Phi_t^1(k_1) \left( \int_{S_t - k_1}^{S_t - k_1 + B_t^*} \beta_t^2 d\Phi_t^2(k_2) \right) + \int_0^{S_t} J'_{t+1}(S_t - k) d\Phi_t^{1,2}(k).
\end{aligned}$$

The second derivative is:

$$\begin{aligned}
G_t^{B''}(S_t) &= (J'_t(-B_t^*) - (p_t^2 + \ell_t^2 - \beta_t^2)) \int_0^{S_t} \phi_t^2(S_t + B_t^* - k_1) d\Phi_t^1(k_1) \\
&\quad - \phi_t^1(S_t) p_t^1 + \phi_t^1(S_t) \int_{B_t^*}^{\infty} p_t^2 d\Phi_t^2(k) - h_t \phi_t^{1,2}(S_t) \\
&\quad - \phi_t^1(S_t) \ell_t^1 + \phi_t^1(S_t) \int_{B_t^*}^{\infty} \ell_t^2 d\Phi_t^2(k) + \int_0^{S_t} J''_{t+1}(S_t - k) d\Phi_t^{1,2}(k) \\
&\quad + \phi_t^1(S_t) \int_0^{B_t^*} \beta_t^2 d\Phi_t^2(k) - \int_0^{S_t} \beta_t^2 \phi_t^2(S_t - k_1) d\Phi_t^1(k_1) \\
&\quad + \phi_t^1(S_t) \int_0^{B_t^*} J'_{t+1}(-k) d\Phi_t^2(k) \\
&\quad + \int_0^{S_t} d\Phi_t^1(k_1) \int_{S_t - k_1}^{S_t + B_t^* - k_1} J''_{t+1}(S_t - k_1 - k_2) d\Phi_t^2(k_2).
\end{aligned}$$

Since  $J''_{t+1} \leq 0$ , and  $p_t^1 > p_t^2$ , and  $\ell_t^1 > \ell_t^2$ , by omitting some terms;



$$\begin{aligned}
G_t^{B''}(S_t) &\leq (J'_t(-B_t^*) - (p_t^2 + \ell_t^2 - \beta_t^2)) \int_0^{S_t} \phi_t^2(S_t + B_t^* - k_1) d\Phi_t^1(k_1) \\
&\quad - \phi_t^1(S_t) \int_0^{B_t^*} p_t^2 d\Phi_t^2(k) - h_t \phi_t^{1,2}(S_t) + \phi_t^1(S_t) \int_0^{B_t^*} \beta_t^2 d\Phi_t^2(k) - \beta_t^2 \phi_t^{1,2}(S_t) \\
&\quad - \phi_t^1(S_t) \int_0^{B_t^*} \ell_t^2 d\Phi_t^2(k) + \phi_t^1(S_t) \int_0^{B_t^*} J'_{t+1}(-k) d\Phi_t^2(k). \\
&= (J'_t(-B_t^*) - (p_t^2 + \ell_t^2 - \beta_t^2)) \int_0^{S_t} \phi_t^2(S_t + B_t^* - k_1) d\Phi_t^1(k_1) \\
&\quad + \phi_t^1(S_t) \int_0^{B_t^*} (J'_t(-k) - (p_t^2 + \ell_t^2 - \beta_t^2)) d\Phi_t^2(k) - (h_t + \beta_t^2) \phi_t^{1,2}(S_t) \leq 0
\end{aligned}$$

Since  $p_t^2 + \ell_t^2 - \beta_t^2 = J'_{t+1}(-B_t^*)$ , and  $p_t^2 + \ell_t^2 - \beta_t^2 > J'_{t+1}(-k)$  where  $k \in [0, B_t^*)$  due to the  $B_t^*$  decision.

Case II:  $B_t^* > q_{t+1}$

Replacing  $B_t^*$  by  $q_{t+1}$  in case 1 and noting that  $p_t^2 + \ell_t^2 - \beta_t^2 > J'_{t+1}(-q_{t+1})$  is enough to conclude that  $G_t^B(S_t)$  is also concave in this case.

7. Let us prove the concavity of  $G_t(S_t)$ .

In each period, we must be in one of the following cases, which are independent of the  $S_t$  values:

- If  $p_t^2 + \ell_t^2 + h_t \geq J'_{t+1}(0) \geq p_t^2 + \ell_t^2 - \beta_t^2$ , we have  $R_t^{1*} = 0$ ,  $R_t^{2*} = 0$ , and  $B_t^* = 0$ , therefore  $R_t^{1,c} = 0$ ,  $R_t^{2,c} = 0$  and  $B_t^c = 0$ ; thus, we have  $G_t(S_t) = G_t^R(S_t) = G_t^B(S_t)$ .
- If  $J'_{t+1}(0) > p_t^2 + \ell_t^2 + h_t$  we have  $R_t^{1*} \geq 0$ ,  $R_t^{2*} \geq 0$ , and  $B_t^* = 0$ , therefore  $R_t^{1,c} \geq 0$ ,  $R_t^{2,c} \geq 0$  and  $B_t^c = 0$ ; thus, we have  $G_t(S_t) = G_t^R(S_t)$ .
- If  $p_t^2 + \ell_t^2 - \beta_t^2 > J'_{t+1}(0)$  we have  $B_t^* \geq 0$ ,  $R_t^{1*} = 0$ , and  $R_t^{2*} = 0$ , therefore  $B_t^c \geq 0$ ,  $R_t^{1,c} = 0$ , and  $R_t^{2,c} = 0$ ; thus, we have  $G_t(S_t) = G_t^B(S_t)$ .

We see that in each period,  $G_t(S_t)$  reduces to some function that is proved to be concave. Therefore  $G_t(S_t)$  is concave.

□

### For the Common Service Strategy

In the proof below, the *CSS* superscript is omitted from the expected profit functions to increase readability.

*Proof.* Let  $j_t(I_t, S_t) = -c_t(S_t - I_t) + G_t(S_t)$ , so  $J_t(I_t) = \max_{S_t: I_t \leq S_t \leq I_t + q_t} j_t(I_t, S_t)$ . We prove by induction.

1. For period  $t = T$ , we have  $B_T = 0$ ,  $R_T = 0$  and  $J_{T+1}(I_T) = v \cdot I_T$ .

$G_T(S_T)$  is concave in  $S_T$ , since  $G'_T(S_T)$  is non-increasing in  $S_T$ :

$$G''_T(S_T) = (v - h_T - p_T^2 - \ell_T)\phi_T^{1,2}(S_T) \leq 0, \text{ since } v < p_T^2.$$

2. Given  $t + 1 \leq T$ , assume that  $G_{t+1}(S_{t+1})$  is concave in  $S_{t+1}$ , then it is easy to see that  $j_{t+1}(I_{t+1}, S_{t+1})$  is jointly concave in  $I_{t+1}$  and  $S_{t+1}$ . So by Lemma B.1,  $J_{t+1}(I_{t+1})$  is concave in  $I_{t+1}$ , and as a result  $J'_{t+1}(I_{t+1})$  is non-increasing in  $I_{t+1}$ .

3. Next let us prove that  $g_t^{CSS}(S_t, R_t, 0)$  is quasi-concave in  $R_t$ . We have,

$$\frac{\partial g_t^{CSS}(S_t, R_t, 0)}{\partial R_t} = (-p_t^2 - \ell_t - h_t + J'_{t+1}(R_t))(1 - \Phi_t^{1,2}(S_t - R_t)) \text{ if } S_t \geq R_t (= 0 \text{ o.w.})$$

If  $R_t^*$  is defined as in (9), we have  $g_t'^{CSS}(S_t, R_t, 0) \geq 0$  when  $0 \leq R_t \leq R_t^*$ , and  $g_t'^{CSS}(S_t, R_t, 0) \leq 0$  when  $R_t > R_t^*$ ; thus,  $g_t^{CSS}(S_t, R_t, 0)$  is quasi-concave with respect to  $R_t$ .  $R_t^*$  is the unique unconstrained optimizer of  $g_t^{CSS}(S_t, R_t, 0)$ , and it is independent of inventory level  $S_t$ .  $R_t^c = \min(R_t^*, S_t)$  maximizes  $g_t^{CSS}(S_t, R_t, 0)$ , for  $0 \leq R_t \leq (S_t)^+$ .

4. Next let us prove that  $g_t^{CSS}(S_t, 0, B_t)$  is quasi-concave in  $B_t$ . Taking the derivative,

$$\frac{\partial g_t^{CSS}(S_t, 0, B_t)}{\partial B_t} = \int_{S_t + B_t / \alpha_t}^{\infty} [p_t^2 - \beta_t^2 + \ell_t^2 - J'_{t+1}(-B_t)] d\Phi_t^{1,2}(D_t^{1,2}).$$

Let us define  $B_t^*$  as in (9), then we have  $g_t'^{CSS}(S_t, 0, B_t) \geq 0$  when  $0 \leq B_t \leq B_t^*$ , and  $g_t'^{CSS}(S_t, 0, B_t) \leq 0$  when  $B_t > B_t^*$ ; thus,  $g_t^{CSS}(S_t, 0, B_t)$  is quasi-concave with respect to  $B_t$ .  $B_t^*$  is the unique unconstrained optimizer of  $g_t^{CSS}(S_t, 0, B_t)$ , and it is independent of inventory level  $S_t$ .  $B_t^c = \min(B_t^*, q_{t+1})$  maximizes  $g_t^{CSS}(S_t, 0, B_t)$ , for  $0 \leq B_t \leq q_{t+1}$ .

5. Let us prove the concavity of  $G_t^R(S_t)$  with respect to  $S_t$ .

We consider  $G_t''^R(S_t)$  in three cases:

(a) Case 1:  $S_t < R_t^*$ :

$G_t''^R(S_t) = J_{t+1}''(S_t) \leq 0$  due to the concavity of  $J_{t+1}$ .

(b) Case 2:  $S_t > R_t^*$ :

$$G_t''^R(S_t) = \int_0^{S_t - R_t^*} J_{t+1}''(S_t - k) d\Phi_t^{1,2}(k) + (J_{t+1}'(R_t^*) - p_t^2 - \ell_t - h_t) \phi_t^{1,2}(S_t - R_t^*) \leq 0$$

due to the choice of  $R_t^*$  and the concavity of  $J_{t+1}$ .

(c) Case 3:  $S_t = R_t^*$  :

$$G_t'^R(R_t^*+) - G_t'^R(R_t^*-) = p_t^2 + \ell_t + h_t - J_{t+1}'(R_t^*) \leq 0$$

due to the choice of  $R_t^*$  and the concavity of  $J_{t+1}$ .

Since  $G_t''^R(S_t) \leq 0$  for all  $S_t$ ,  $G_t^R(S_t)$  is concave in  $S_t$ .

6. Let us prove the concavity of  $G_t^B(S_t)$  with respect to  $S_t$ , where  $G_t^B(S_t) = g_t^{CSS}(S_t, 0, B_t^c)$ .

We have,

$$\begin{aligned} G_t''^B(S_t) &= -\alpha_t(p_t^2 + \ell_t^2 - \beta_t^2 - J_{t+1}'(-B_t)) \phi_t^{1,2}(S_t + B_t/\alpha_t) \\ &\quad + [\alpha_t(p_t^2 + \ell_t^1 - \beta_t^2 - J_{t+1}'(0)) - (p_t^2 + \ell_t^1 + h_t - J_{t+1}'(0))] \phi_t^{1,2}(S_t) \\ &\quad + \int_0^{S_t} J_{t+1}''(S_t - k) d\Phi_t^{1,2}(k) + \int_{S_t}^{S_t + B_t/\alpha_t} \alpha_t^2 J_{t+1}''(\alpha_t(S_t - k)) d\Phi_t^{1,2}(k). \end{aligned}$$

The first term in  $G_t''^B(S_t)$  is negative due to the choice of  $B_t^*$ . The third and the fourth terms in  $G_t''^B(S_t)$  are negative due to the concavity of  $J_{t+1}(S_t)$ . We have  $G_t''^B(S_t) \leq 0$  and therefore,  $G_t^B(S_t, B_t)$  is concave in  $S_t$ .

7. Let us prove the concavity of  $G_t(S_t)$ . In each period, we must be in one of the following cases, which are independent of the  $S_t$  values:

- If  $p_t^2 + \ell_t^2 - \beta_t^2 \leq J_{t+1}'(0) \leq p_t^2 + \ell_t + h_t$ , we have  $R_t^* = B_t^* = 0$ , therefore  $R_t^c = B_t^c = 0$ ; thus, we have  $G_t(S_t) = G_t^R(S_t) = G_t^B(S_t)$ .
- If  $J_{t+1}'(0) > p_t^2 + \ell_t + h_t$ , we have  $R_t^* \geq 0$  and  $B_t^* = 0$ , therefore  $R_t^c \geq 0$  and  $B_t^c = 0$ ; thus, we have  $G_t(S_t) = G_t^R(S_t) \geq G_t^B(S_t)$ .
- If  $p_t^2 + \ell_t^2 - \beta_t^2 < J_{t+1}'(0)$ , we have  $B_t^* \geq 0$  and  $R_t^* = 0$ , therefore  $B_t^c \geq 0$  and  $R_t^c = 0$ ; thus, we have  $G_t(S_t) = G_t^B(S_t) \geq G_t^R(S_t)$ .

We see that in each period,  $G_t(S_t)$  reduces to some function that is proved to be concave. Therefore  $G_t(S_t)$  is concave.

□

## APPENDIX C

### PROOFS FOR CHAPTER 4

#### Proofs for Infinite Capacity Case

##### Proof of Theorem 4.1

*Proof.* (i) In order to find the optimal leadtime, we look at the first order condition. Taking the derivative of (11) with respect to  $\ell$  and rearranging terms, we get the first order condition as:

$$\frac{\partial \Pi}{\partial \ell} = \frac{1}{\mu} e^{-(\alpha+\mu)\ell} (B - Ae^{\mu\ell} - Ce^{-\mu\ell}) = 0,$$

where  $A = \alpha\mu R$ ,  $B = (\alpha + \mu)(c + \mu R)$ , and  $C = (\alpha + 2\mu)c$ . By letting  $x = e^{\mu\ell}$ , it is easily seen that this is a quadratic equation. So, the leadtime values that satisfy the first order condition are provided below;

$$\begin{aligned} \ell^{\pm} &= \frac{1}{\mu} \ln(x^{\pm}) = \frac{1}{\mu} \ln \frac{B \pm \sqrt{B^2 - 4AC}}{2A} \\ &= \frac{1}{\mu} \ln \left( \frac{(\alpha + \mu)(c + \mu R) \pm \sqrt{(\alpha + \mu)^2(c - \mu R)^2 + 4\mu^3 Rc}}{2\alpha R \mu} \right). \end{aligned}$$

$(\alpha + \mu)^2(c - \mu R)^2 + 4\mu^3 Rc > 0$  holds for any positive  $\alpha$ ,  $\mu$ ,  $c$  and  $R$ , hence two real roots exist.

Profit function is given by:

$$\Pi(\ell) = e^{-\alpha\ell} \left( R - \frac{c}{\mu} e^{-\mu\ell} \right) (1 - e^{-\mu\ell}).$$

When  $\ell = 0$ , the arrival rate  $(1 - e^{-\mu\ell}) = 0$ , therefore, we have  $\Pi(0) = 0$ . Since  $\ell$ ,  $\alpha$  and  $\mu$  are greater than zero, the only part of the  $\Pi(\ell)$  that needs to be considered for positive profit is the middle part. Below, the requirement (23) gives the condition on  $\ell$  that ensures having positive profit.

$$\Pi(\ell) > 0 \Leftrightarrow \left( R - \frac{c}{\mu} e^{-\mu\ell} \right) > 0 \Leftrightarrow \ell > \ln\left(\frac{c}{\mu R}\right) \frac{1}{\mu} \quad (23)$$

By omitting some positive parts from  $\ell^+$ , we can show that  $\ell^+$  satisfies (23), and  $\Pi(\ell^+) > 0$ :

$$\ell^+ > \frac{1}{\mu} \ln \left( \frac{(\alpha + \mu)(c + \mu R) + \sqrt{(\alpha + \mu)^2(c - \mu R)^2}}{2\alpha R\mu} \right) = \frac{1}{\mu} \ln \left( \frac{c(\alpha + \mu)}{\alpha R\mu} \right) > \left( \ln \frac{c}{\mu R} \right) \frac{1}{\mu}.$$

Note that, when  $\ell = 0$ ,  $\frac{\partial \Pi}{\partial \ell} = \frac{B-A-C}{\mu} = \mu R - c$ . Therefore, when  $\mu R < c$ , the profit function  $\Pi(\ell)$  is decreasing at zero, which ensures the small root  $\ell^-$  is the minimizer and the maximizer is the bigger root  $\ell^+$  due to the fact that  $\Pi(\ell^+) > 0$ . When  $\mu R \geq c$ , we will show that the small root  $\ell^-$  is smaller than or equal to zero, therefore we have only one positive root  $\ell^+$  and it is the maximizer. For this result, we need to show that  $x^- \leq 1$ .

Observe that  $\sqrt{(\alpha + \mu)^2(c - \mu R)^2 + 4\mu^3 R c} > \sqrt{\mu^2(c + \mu R)^2}$  by omitting the positive term  $(\alpha^2 + 2\alpha\mu)(c - \mu R)^2$  from the left hand side. Therefore;

$$\begin{aligned} x^- &= \frac{(\alpha + \mu)(c + \mu R) - \sqrt{(\alpha + \mu)^2(c - \mu R)^2 + 4\mu^3 R c}}{2\alpha R\mu} \\ &< \frac{(\alpha + \mu)(c + \mu R) - \sqrt{\mu^2(c + \mu R)^2}}{2\alpha R\mu} = \frac{(c + \mu R)}{2R\mu} \leq 1 \end{aligned}$$

This concludes the proof of  $\ell^+$  being the maximizer of (11).

(ii) Let  $R = \frac{cn}{\mu}$  where  $n > 0$ . Then we have,

$$\theta = \frac{\alpha + \alpha n + \mu + \mu n + \sqrt{(\alpha^2 + 2\alpha\mu)(n-1)^2 + \mu^2(1+n)^2}}{2\alpha n}$$

$\ell^* = \frac{1}{\mu} \ln(\theta)$  is increasing in  $\mu$ , and since  $R = \frac{cn}{\mu}$ ,  $\ell^*$  is decreasing in  $R$ . Since  $\lim_{R \rightarrow \infty} \theta = \frac{\alpha + \mu}{\alpha} > 1$ ,  $\ell^*$  is always positive.

(iii) We show this by omitting some positive terms from  $\ell^*$ .

$$\begin{aligned} \ell^* &= \frac{1}{\mu} \left( \ln \frac{(\alpha + \mu)(c + \mu R) + \sqrt{(\alpha + \mu)^2(c - \mu R)^2 + 4\mu^3 R c}}{2\alpha R\mu} \right)^+ \\ &> \frac{1}{\mu} \left( \ln \frac{(\alpha + \mu)(c + \mu R) + \sqrt{(\alpha + \mu)^2(c - \mu R)^2}}{2\alpha R\mu} \right)^+ \\ &= \frac{1}{\mu} \left( \ln \frac{(\alpha + \mu)c}{\alpha R\mu} \right)^+ = \ell^N. \end{aligned}$$

□

## Proofs for Finite Capacity Case

The following three lemmas are used for the proofs of Theorem 4.2 and Theorem 4.3.

**Lemma C.1.** Let  $E_{cost}(i) = c \cdot i + E_{clearing}(i)$ , for  $i = 0, 1, \dots$ , where,

$$E_{clearing}(i) = \frac{Pr\{X(\ell) = 0\} \cdot c \cdot i + \sum_{j=1}^i Pr\{X(\ell) = j\}[(i-j) \cdot c + E_{clearing}(i-j)]}{1 - Pr\{X(\ell) = 0\}}.$$

Then,  $E_{clearing}(i)$  and  $E_{cost}(i)$  are increasing and convex in  $i$ .

*Proof.* Proof will be done by induction on  $i$  for  $E_{clearing}(i)$  and this will imply that  $E_{cost}(i)$  is also increasing and convex in  $i$  since  $c \cdot i$  is convex and increasing in  $i$ . Let  $\Delta E_{clearing}(i) = E_{clearing}(i+1) - E_{clearing}(i)$ . For the initial step,  $E_{clearing}(1) = \frac{c \cdot Pr\{X(\ell)=0\}}{1 - Pr\{X(\ell)=0\}} > E_{clearing}(0) = 0$ , and  $\Delta E_{clearing}(1) > \Delta E_{clearing}(0)$  which is given by;

$$\frac{c \cdot Pr\{X(\ell) = 0\} + Pr\{X(\ell) = 1\}[c + E_{clearing}(1)]}{1 - Pr\{X(\ell) = 0\}} > \frac{c \cdot Pr\{X(\ell) = 0\}}{1 - Pr\{X(\ell) = 0\}}.$$

Assume  $E_{clearing}(i)$  is increasing and convex in  $i$  for  $i = 1, \dots, (n-1)$ . Given the general term for  $\Delta E_{clearing}(i)$  below,

$$\frac{Pr\{X(\ell) = 0\} \cdot c + \sum_{j=1}^i Pr\{X(\ell) = j\}[c + E_{clearing}(i+1-j) - E_{clearing}(i-j)]}{1 - Pr\{X(\ell) = 0\}},$$

it is easily seen that  $E_{clearing}(n) > E_{clearing}(n-1)$  and  $\Delta E_{clearing}(n) > \Delta E_{clearing}(n-1)$  due to the induction hypothesis.  $\square$

**Lemma C.2 (from [76]).** Let  $\phi(i) = g(f(i))$ ,  $i = 0, 1, \dots$ , where  $f(i)$  is a convex, non-decreasing, integer-valued function of  $i = 0, 1, \dots$ , with  $f(0) \geq 0$ , and  $g(j)$  is a concave, non-increasing function of  $j = 0, 1, \dots$ . Then  $\phi(i)$  is a concave non-increasing function of  $i = 0, 1, \dots$

**Lemma C.3 (from [76]).** Let  $g_k(i) = \max_{a=0,1,\dots,k}\{ar + f(i+a)\}$ ,  $i = 0, 1, \dots$ . If  $f(\cdot)$  is concave and non-increasing, then  $g_k(\cdot)$  is concave and non-increasing in  $i$ .

These three lemmas are used for the proof of Theorem 4.3.

**Lemma C.4 (from [81]).** The collection of distribution functions  $F_s(k)$  is stochastically increasing in  $s$  on  $S$ , if and only if  $\int h(k)dF_s(k)$  is increasing in  $s$  on  $S$  for each increasing real-valued function  $h(k)$ .

**Lemma C.5.** Given  $\{F_s(k) : s \in S\}$  is a collection of stochastically increasing distribution functions in  $s$ , and  $v(i, k, s)$  is a real-valued function that is non-decreasing in  $s$  for  $i$  and some fixed  $k$ , and also non-decreasing in  $k$  for  $i$  and some fixed  $s$  then;

$$V(i, s) = \int v(i, k, s) dF_s(k) \text{ is non-decreasing in } s.$$

*Proof.* We have,

$$\begin{aligned} V(i, s+1) &= \int v(i, k, s+1) dF_{s+1}(k) \geq \int v(i, k, s) dF_{s+1}(k) \\ &\geq \int v(i, k, s) dF_s(k) = V(i, s) \end{aligned}$$

The first inequality follows from  $v(i, k, s)$  being non-decreasing in  $s$ , and the second inequality follows from Lemma C.4 since  $v(i, k, s)$  is non-decreasing in  $k$  and  $F_s(k)$  is stochastically increasing in  $s$ .  $\square$

**Lemma C.6.** Let  $v_{n-1}(i, k, s) = \max_{a=0, \dots, k} \{a \cdot R + U_{n-1}(i+a, s)\}$ . Assuming  $U_{n-1}(i, s+1) - U_{n-1}(i, s) \geq 0$ ,  $v_{n-1}(i, k, s)$  is non-decreasing in  $k$  and  $s$ .

*Proof.*  $U_{n-1}(i, s+1) \geq U_{n-1}(i, s) \Rightarrow a \cdot R + U_{n-1}(i+a, s+1) \geq a \cdot R + U_{n-1}(i+a, s)$  for every  $a$ . Then the maximization for  $s+1$  is applied over a set whose elements have a higher value than that of  $s$ , which implies  $v_{n-1}(i, k, s)$  is non-decreasing in  $s$ . Obviously,  $v_{n-1}(i, k, s)$  is also non-decreasing in  $k$ , since the feasible set gets larger as  $k$  increases.  $\square$

### Proof of Theorem 4.2

*Proof.* The proof will be done by induction on  $n$ . The case  $n = 1$  is immediate, since  $V_0(i) = -E_{cost}(i)$  is concave and non-increasing in  $i$  by Lemma C.1. Taking expectations and using Lemma C.2,  $U_1(i)$  is concave and non-increasing in  $i$ , since  $(i - X(\ell))^+$  is convex and non-decreasing in  $i$ , for each fixed value of  $X(\ell)$ . Then using lemma C.3 and taking expectations, we get  $V_1(i)$  as concave and non-increasing in  $i$ .

Now, suppose  $n \geq 2$  and  $V_{n-1}(i)$  is concave and non-increasing in  $i$ .

Since  $(i - X(\ell))^+$  is convex and non-decreasing in  $i$ , for each fixed value of  $X(\ell)$ , it follows from the induction hypothesis, using Lemma C.2 and taking expectations that  $U_n(i)$  is concave and non-increasing in  $i$ . So, again by using Lemma C.3 and taking expectations, we get  $V_n(i)$  as concave and non-increasing in  $i$ , which concludes the proof.  $\square$

### Proof of Theorem 4.3

*Proof.* First, we show that the optimal acceptance policy is a threshold policy in  $i$  for a fixed service index  $s$ , and in the second part, we prove that the optimal decision levels are convex in  $s$  for a fixed  $i$ , for any period  $n$ .

1. Let us start by showing that the optimal acceptance policy is a threshold policy in  $i$  for a fixed service index  $s$ . We first need to prove the following lemma:

**Lemma C.7.** *Let  $g(i, s) = E\left\{f(0, s^+)I\{X(\ell) \geq i\} + f((i - X(\ell))^+, s^-)I\{X(\ell) < i\}\right\}$ ,  $i = 0, 1, \dots$ , where  $f(i, s)$  is concave and non-increasing in  $i$  for a fixed  $s$ . Then,  $g(i, s)$  is concave and non-increasing in  $i$ ,  $i = 0, 1, \dots$ , for a fixed  $s$ .*

*Proof.*  $g(i, s) = f(0, s^+) \sum_{j=i}^{\infty} Pr\{X(\ell) = j\} + \sum_{j=0}^{i-1} f(i - j, s^-) Pr\{X(\ell) = j\}$

If we define  $\Delta g(i, s) = g(i, s) - g(i + 1, s)$  and  $\Delta f(i, s^-) = f(i, s^-) - f(i + 1, s^-)$ , we have;

$$\Delta g(i, s) = (f(0, s^+) - f(1, s^-))Pr\{X(\ell) = i\} + \sum_{j=0}^{i-1} \Delta f(i - j, s^-)Pr\{X(\ell) = j\}.$$

For concavity, we need  $\Delta g(i, s) \leq \Delta g(i + 1, s)$  for every  $i = 0, 1, \dots$

$$\begin{aligned} \Delta g(i + 1, s) &= (f(0, s^+) - f(1, s^-))Pr\{X(\ell) = i + 1\} \\ &\quad + \sum_{j=0}^i \Delta f(i + 1 - j, s^-)Pr\{X(\ell) = j\} \\ &\geq (f(0, s^+) - f(1, s^-))Pr\{X(\ell) = i + 1\} \\ &\quad + \sum_{j=0}^{i-1} \Delta f(i + 1 - j, s^-)Pr\{X(\ell) = j\} \\ &\geq (f(0, s^+) - f(1, s^-))Pr\{X(\ell) = i\} + \sum_{j=0}^{i-1} \Delta f(i - j, s^-)Pr\{X(\ell) = j\} \\ &= \Delta g(i, s) \end{aligned}$$

First inequality is due to non-increasingness of  $f(i, s)$  for a fixed  $s$ . Second inequality comes from two properties. First one is the concavity of the  $f(i, s)$ , which ensures  $\Delta f(i + 1 - j, s^-) \geq \Delta f(i - j, s^-)$ . Second one is due to the batch processing assumption, which ensures  $Pr\{X(\ell) = i\} = Pr\{X(\ell) = i + 1\}$  where  $i \leq M - 1$  due to assumption of not accepting more than  $M$  customers to the system.  $\square$

Showing that the optimal acceptance policy is a threshold policy in  $i$  for a fixed service index  $s$  is very similar to the proof of Theorem 4.2. The only difference is instead of



Lemma C.2, Lemma C.7 is used when the service index is fixed.

2. The proof is done by mathematical induction on  $n$ . Let  $h_n(i, s) = U_n(i, s+1) - U_n(i, s)$ .

We show that:

- $h_n(i, s) \geq 0, \forall i, s$  ( $\star$ ),
- $h_n(i, 0)$  is decreasing and  $h_n(i, 1)$  is increasing in  $i$  ( $\star\star$ ),

holds for any period  $n$ , which implies that the optimal decision levels are convex in  $s$  for any period  $n$ . First we start by proving a result that shows how the structure of  $h_n(i, 0)$  and  $h_n(i, 1)$  determines the form of the acceptance policy, and then proceed with the proof.

**Lemma C.8.** *Given that  $h_n(i, 0)$  is decreasing in  $i$ , and  $h_n(i, 1)$  is increasing in  $i$ , the optimal decisions in period  $n$  have the property of being convex in  $s$ ;  $a_{s=0}^*(i) \geq a_{s=1}^*(i)$  and  $a_{s=1}^*(i) \leq a_{s=2}^*(i)$ .*

*Proof.* Let us start with  $h_n(i, 0)$ .

$$\begin{aligned} h_n(i, 0) &= U_n(i, 1) - U_n(i, 0) \geq U_n(i+1, 1) - U_n(i+1, 0) = h_n(i+1, 0) & \forall i \\ U_n(i, 1) - U_n(i+1, 1) &\geq U_n(i, 0) - U_n(i+1, 0) & \forall i \end{aligned} \quad (24)$$

We know from the first part of this theorem that  $U_n(i, \cdot)$  is concave and non-increasing in  $i$ . Then the optimal acceptance decision for a service level  $s$  is decided by the intersection of  $R$  and the  $U_n(i, s) - U_n(i+1, s)$  curve. Inequality 24 shows that  $a_{s=0}^*(i) \geq a_{s=1}^*(i)$ , since  $U_n(i, 1) - U_n(i+1, 1)$  intersects  $R$  first. The proof of the case for  $a_{s=1}^*(i) \leq a_{s=2}^*(i)$  is very similar to the first case.  $\square$

Since we have only three service levels:

$$h_n(i, s) = \begin{cases} [V_{n-1}(0, 2) - V_{n-1}(0, 1)]Pr\{X(\ell) \geq i\} & \text{if } s = 0 \\ \sum_{j=0}^{i-1} (V_{n-1}(i-j, 1) - V_{n-1}(i-j, 0))Pr\{X(\ell) = j\} & \text{if } s = 1 \end{cases}$$

We have  $V_0(i, s+1) - V_0(i, s) = S_{last}$  by assumption. So  $h_1(i, s) \geq 0 \forall i, s$  for period

1. Also, it is easily seen that  $h_1(i, 0)$  is decreasing in  $i$ , and  $h_1(i, 1)$  is increasing in  $i$ .

By Lemma C.8,  $a_{s=0}^*(i) \geq a_{s=1}^*(i)$  and  $a_{s=1}^*(i) \leq a_{s=2}^*(i)$  for period 1.

Assume  $(\star)$  and  $(\star\star)$  hold for period  $(n - 1)$ .

By Lemma C.6,  $v_{n-1}(i, k, s) = \max_{a=0,\dots,k} \{a \cdot R + U_{n-1}(i + a, s)\}$  is non-decreasing in  $k$  and  $s$ . Then, by Lemma C.5,  $V_{n-1}(i, s)$  is non-decreasing in  $s$ . So for period  $n$ , since  $V_{n-1}(i, s+1) \geq V_{n-1}(i, s)$ , it is easily seen that  $h_n(i, s) = U_n(i, s+1) - U_n(i, s) \geq 0$   $(\star)$  and  $h_n(i, 0) = U_n(i, 1) - U_n(i, 0)$  is decreasing in  $i$ . To see that  $h_n(i, 1)$  is increasing in  $i$ , let us look at  $h_n(0, 1)$  for the sequence of  $i$ 's:

$$\begin{aligned}
h_n(0, 1) &= 0 \\
h_n(1, 1) &= [V_{n-1}(1, 1) - V_{n-1}(1, 0)]Pr\{X(\ell) = 0\} \\
h_n(2, 1) &= [V_{n-1}(2, 1) - V_{n-1}(2, 0)]Pr\{X(\ell) = 0\} \\
&\quad + [V_{n-1}(1, 1) - V_{n-1}(1, 0)]Pr\{X(\ell) = 1\} \\
h_n(3, 1) &= [V_{n-1}(3, 1) - V_{n-1}(3, 0)]Pr\{X(\ell) = 0\} \\
&\quad + [V_{n-1}(2, 1) - V_{n-1}(2, 0)]Pr\{X(\ell) = 1\} \\
&\quad + [V_{n-1}(1, 1) - V_{n-1}(1, 0)]Pr\{X(\ell) = 2\} \\
&\vdots
\end{aligned}$$

$h_n(i, 1)$  is increasing in  $i$ , since  $Pr\{X(\ell) = i - 1\} = Pr\{X(\ell) = i\}$ . Using Lemma C.8, we conclude that  $a_{s=0}^*(i) \geq a_{s=1}^*(i)$  and  $a_{s=1}^*(i) \leq a_{s=2}^*(i)$  for period  $n$ .

□

## APPENDIX D

### PROOFS FOR CHAPTER 5

#### Proof of Theorem 5.1

*Proof.* Assume that there exists a function satisfying the conditions in the theorem. We will show that  $\bar{V}$  is equal to  $\tilde{V}$ . Let us use the same kind of martingales as in expressions (14) and (15) for  $\bar{V}(t, n(t))$ . For  $s \geq t$ , let:

$$m(s) = \bar{V}(s, [n(t) - N_B(s) + N_B(t)]^+) - \bar{V}(t, n(t)) - \int_t^s \mathcal{G}\bar{V}(u, [n(t) - N_B(u) + N_B(t)]^+) du.$$

$m(s)$  is a martingale by Dynkin's Lemma and since the expected value of this martingale at any time  $s$  is equal to its expected value at the starting time  $t$ , we have  $Em(s) = 0$ . Further, by the optional sampling theorem, for any stopping time  $\tau \geq t$  we have:

$$E[\bar{V}(\tau, [n(t) - N_B(\tau) + N_B(t)]^+)] - E \int_t^\tau \mathcal{G}\bar{V}(u, [n(t) - N_B(u) + N_B(t)]^+) du = \bar{V}(t, n(t)) \quad (25)$$

$$\begin{aligned} & E[\bar{V}(\tau, [n(t) - N_B(\tau) + N_B(t)]^+)] \quad (26) \\ & - E \int_t^\tau [\mathcal{G}(\bar{V} + \Pi)(u, [n(t) - N_B(u) + N_B(t)]^+) + \lambda_{BPB} \mathbf{1}_{\{N_B(u) - N_B(t) < n(t)\}}] du \\ & = \bar{V}(t, n(t)) - E \int_t^\tau [\mathcal{G}\Pi(u, [n(t) - N_B(u) + N_B(t)]^+) + \lambda_{BPB} \mathbf{1}_{\{N_B(u) - N_B(t) < n(t)\}}] du. \end{aligned}$$

If we subtract  $E \int_t^\tau [\mathcal{G}\Pi(u, [n(t) - N_B(u) + N_B(t)]^+) + \lambda_{BPB} \mathbf{1}_{\{N_B(u) - N_B(t) < n(t)\}}] du$  from both sides of (25), the left-hand side of the resulting term, given by (26), is always positive by conditions (i), (iii) and (iv). Therefore,

$$\bar{V}(t, n(t)) \geq E \int_t^\tau [\mathcal{G}\Pi(u, [n(t) - N_B(u) + N_B(t)]^+) + \lambda_{BPB} \mathbf{1}_{\{N_B(u) - N_B(t) < n(t)\}}] du.$$

Since  $\bar{V}(t, n(t))$  is greater than or equal to each term in the right-hand side of equation (19) for any  $\tau$ , it is also greater than or equal to the supremum over all  $\tau$ , which is  $\tilde{V}(t, n(t))$  in equation (19). Hence, we conclude that  $\bar{V}(t, n(t)) \geq \tilde{V}(t, n(t))$  for any stopping time  $\tau \geq t$ .

To prove that  $\bar{V}(t, n(t)) \leq \tilde{V}(t, n(t))$ , we will define a specific stopping time. Let  $\sigma$  be defined as  $\sigma = \inf\{t \leq s \leq T : \bar{V}(s, [n(t) - N_B(s) + N_B(t)]^+) = 0\}$ . Note that  $\sigma$  is well-defined because  $\bar{V}(T, \cdot) = 0$ . Replacing  $\tau$  in equation (26) with the specific stopping time  $\sigma$ , we obtain:

$$\begin{aligned} & E[\bar{V}(\sigma, [n(t) - N_B(\sigma) + N_B(t)]^+)] \\ & - E \int_t^\sigma [\mathcal{G}(\bar{V} + \Pi)(u, [n(t) - N_B(u) + N_B(t)]^+) + \lambda_{BPB} \mathbf{1}_{\{N_B(u) - N_B(t) < n(t)\}}] du \\ & = \bar{V}(t, n(t)) - E \int_t^\sigma [\mathcal{G}\Pi(u, [n(t) - N_B(u) + N_B(t)]^+) + \lambda_{BPB} \mathbf{1}_{\{N_B(u) - N_B(t) < n(t)\}}] du. \end{aligned} \quad (27)$$

The definition of  $\sigma$  implies  $\bar{V}(\sigma, [n(t) - N_B(\sigma) + N_B(t)]^+) = 0$ , and the definition of  $\sigma$  and condition (iv) together imply that  $\mathcal{G}(\Pi)(u, [n(t) - N_B(u) + N_B(t)]^+) + \lambda_{BPB} = 0$  for all  $u \in [t, \sigma]$ . Therefore the left-hand side of (27) is zero and we have:

$$\begin{aligned} \bar{V}(t, n(t)) & = E \int_t^\sigma [\mathcal{G}\Pi(u, [n(t) - N_B(u) + N_B(t)]^+) + \lambda_{BPB} \mathbf{1}_{\{N_B(u) - N_B(t) < n(t)\}}] du \\ & \leq \tilde{V}(t, n(t)). \end{aligned}$$

The inequality follows from the fact that the left-hand side of the inequality is the right-hand side of equation (19) for a specific stopping time, and  $\tilde{V}(t, n(t))$  is the supremum over all stopping times  $\tau$  in that equation. Hence,  $\bar{V}(t, n(t)) = \tilde{V}(t, n(t))$ .  $\square$

### Proof of Lemma 5.1

*Proof.* Using the definition of derivative, conditioning on  $N_i(T) - N_i(t + h)$  and omitting the zero terms, we get:

$$\begin{aligned} \frac{\partial P[N_i(T) - N_i(t) \geq k]}{\partial t} & = - \lim_{h \rightarrow 0} \frac{P[N_i(T) - N_i(t) \geq k] - P[N_i(T) - N_i(t + h) \geq k]}{h} \\ & = - \lim_{h \rightarrow 0} \frac{\lambda_i h P(N_i(T) - N_i(t) = k - 1)}{h} \\ & = -\lambda_i P(N_i(T) - N_i(t) = k - 1). \end{aligned} \quad (28)$$

$$= -\lambda_i P(N_i(T) - N_i(t) = k - 1). \quad (29)$$

Therefore, we have  $\frac{\partial \sum_{k=1}^{n(t)} P[N_i(T) - N_i(t) \geq k]}{\partial t} = -\lambda_i P[N_i(T) - N_i(t) \leq n(t) - 1]$ . Using

this equality we have:

$$\begin{aligned}
\mathcal{G}\Pi(t, n(t)) &= \frac{\partial \Pi(t, n(t))}{\partial t} + \lambda_B [\Pi(t, n(t) - 1) - \Pi(t, n(t))] \\
&= -\lambda_1 p_1 P[N_1(T) - N_1(t) \leq n(t) - 1] - \lambda_2 p_2 P[N_2(T) - N_2(t) \leq n(t) - 1] \\
&\quad - \lambda_B p_1 P[N_1(T) - N_1(t) \geq n(t)] - \lambda_B p_2 P[N_2(T) - N_2(t) \geq n(t)] \\
&= -\lambda_1 p_1 (1 - P[N_1(T) - N_1(t) \geq n(t)]) - \lambda_2 p_2 (1 - P[N_2(T) - N_2(t) \geq n(t)]) \\
&\quad - \lambda_B p_1 P[N_1(T) - N_1(t) \geq n(t)] - \lambda_B p_2 P[N_2(T) - N_2(t) \geq n(t)] \\
&= -\lambda_1 p_1 - \lambda_2 p_2 + p_1 (\lambda_1 - \lambda_B) P[N_1(T) - N_1(t) \geq n(t)] \\
&\quad + p_2 (\lambda_2 - \lambda_B) P[N_2(T) - N_2(t) \geq n(t)].
\end{aligned}$$

□

## Proof of Theorem 5.2

*Proof.* The proof will be done by mathematical induction on  $n(t)$ . When  $n(t) = 1$ ,  $\bar{V}(t, n(t) - 1) = 0$ , and we have  $L(t, 1) = \mathcal{G}\Pi(t, 1) + \lambda_B p_B$ , which is an increasing function in  $t$  since  $\mathcal{G}\Pi(t, 1)$  is an increasing function in  $t$ . We claim that for  $t \leq x_1$ :

$$L(t, 1) = \mathcal{G}\Pi(t, 1) + \lambda_B p_B \leq \mathcal{G}\Pi(x_1, 1) + \lambda_B p_B \leq 0.$$

The first inequality is by the increasing property of  $\mathcal{G}\Pi(t, 1)$  in  $t$ . The second inequality follows from the fact that if  $\mathcal{G}\Pi(x_1, 1) + \lambda_B p_B > 0$  then  $\int_{x_1}^T L(s, n(t)) e^{-\lambda_B(s-t)} ds > 0$ , which contradicts the definition of  $x_1$ . Hence for  $t \leq x_1$  (or  $\bar{V}(t, 1) = 0$  by the definition of  $\bar{V}$ ):

$$\begin{aligned}
\mathcal{G}(\bar{V} + \Pi)(t, 1) + \lambda_B p_B &= \frac{\partial \bar{V}(t, 1)}{\partial t} + \lambda_B [\bar{V}(t, 0) - \bar{V}(t, 1)] + \mathcal{G}\Pi(t, 1) + \lambda_B p_B \\
&= \mathcal{G}\Pi(t, 1) + \lambda_B p_B = L(t, 1) \leq 0.
\end{aligned}$$

Thus, condition (iii) is satisfied when  $n(t) = 1$  and  $t \leq x_1$  (or  $\bar{V}(t, 1) = 0$ ).

When  $t > x_1$  (or  $\bar{V}(t, 1) > 0$ ):

$$\begin{aligned}
\mathcal{G}(\bar{V} + \Pi)(t, 1) + \lambda_B p_B &= \frac{\partial \bar{V}(t, 1)}{\partial t} + \lambda_B [\bar{V}(t, 0) - \bar{V}(t, 1)] + \mathcal{G}\Pi(t, 1) + \lambda_B p_B \\
&= \frac{\partial \bar{V}(t, 1)}{\partial t} - \lambda_B \bar{V}(t, 1) + \mathcal{G}\Pi(t, 1) + \lambda_B p_B.
\end{aligned} \tag{30}$$

By the definition of  $\bar{V}(t, n(t))$ , we have  $\bar{V}(t, 1) = \int_t^T L(s, 1)e^{-\lambda_B(s-t)} ds$ . Taking the derivative with respect to  $t$ , we get:

$$\frac{\partial \bar{V}(t, 1)}{\partial t} = \int_t^T \lambda_B L(s, 1)e^{-\lambda_B(s-t)} ds - L(t, 1) = \lambda_B \bar{V}(t, 1) - \mathcal{G}\Pi(t, 1) - \lambda_B p_B. \quad (31)$$

Substituting (31) into (30), we get  $\mathcal{G}(\bar{V} + \Pi)(t, 1) + \lambda_B p_B = 0$ .  $\bar{V}(t, 1) > 0$ . Therefore condition (iv) is satisfied when  $n(t) = 1$ . Moreover, we have  $\bar{V}(t, 1) \geq \bar{V}(t, 0) = 0$  by the definition of  $x_1$  ( $\exists t: \bar{V}(t, 1) > 0$  if  $x_1 > 0$ ).

Now assume that the following statements hold for  $n(t) \leq k < M$ : there exist  $k$  time thresholds with  $T \geq x_1 \geq \dots \geq x_k \geq 0$  such that  $\bar{V}(t, n(t))$  is derived from equation (22) and satisfies conditions (i)-(iv), and the inequality  $\bar{V}(t, n(t)) \geq \bar{V}(t, n(t) - 1)$  holds for  $n(t) = 1 \dots k$ .

For  $n(t) = k + 1$ ,

$$L(t, k + 1) = \mathcal{G}\Pi(t, k + 1) + \lambda_B p_B + \lambda_B \bar{V}(t, k) \geq \mathcal{G}\Pi(t, k) + \lambda_B p_B + \lambda_B \bar{V}(t, k - 1) = L(t, k),$$

since  $\mathcal{G}\Pi(t, k)$  and  $\bar{V}(t, k)$  are increasing in  $k$  by the induction assumption. This implies:

$$\int_t^T L(s, k + 1)e^{-\lambda_B(s-t)} ds \geq \int_t^T L(s, k)e^{-\lambda_B(s-t)} ds.$$

Together with equation (22), this implies  $\bar{V}(t, k + 1) \geq \bar{V}(t, k)$  and  $x_k \geq x_{k+1}$ .

For  $t \leq x_{k+1}$  (or  $\bar{V}(t, k + 1) = 0$ ),

$$\begin{aligned} \mathcal{G}(\bar{V} + \Pi)(t, k + 1) + \lambda_B p_B &= \frac{\partial \bar{V}(t, k + 1)}{\partial t} + \lambda_B [\bar{V}(t, k) - \bar{V}(t, k + 1)] + \mathcal{G}\Pi(t, k + 1) + \lambda_B p_B \\ &= \mathcal{G}\Pi(t, k + 1) + \lambda_B p_B + \lambda_B \bar{V}(t, k) = L(t, k + 1) \\ &\leq L(x_{k+1}, k + 1) \leq 0. \end{aligned}$$

Note that  $\bar{V}(t, k) = \bar{V}(t, k + 1) = 0$  since  $t \leq x_{k+1} \leq x_k$ . The first inequality follows from  $\mathcal{G}\Pi(t, k + 1)$  being increasing in  $t$ , and the second inequality follows from the fact that if  $L(x_{k+1}, k + 1) > 0$  then this will contradict the definition of  $x_{k+1}$ . Therefore, condition (iii) is satisfied, when  $t \leq x_{k+1}$  (or  $\bar{V}(t, k + 1) = 0$ ).

For  $t > x_{k+1}$  (or  $\bar{V}(t, k+1) > 0$ ),

$$\begin{aligned}
& \mathcal{G}(\bar{V} + \Pi)(t, k+1) + \lambda_B p_B \\
&= \frac{\partial \bar{V}(t, k+1)}{\partial t} + \lambda_B [\bar{V}(t, k) - \bar{V}(t, k+1)] + \mathcal{G}\Pi(t, k+1) + \lambda_B p_B \\
&= -L(t, k+1) + \lambda_B \bar{V}(t, k+1) + \lambda_B [\bar{V}(t, k) - \bar{V}(t, k+1)] + \mathcal{G}\Pi(t, k+1) + \lambda_B p_B = 0.
\end{aligned}$$

Therefore condition (iv) is satisfied when  $t > x_{k+1}$  (or  $\bar{V}(t, k+1) = 0$ ).

For  $n(t) = k+1$  we showed that conditions (i)-(iv) hold. Thus  $\bar{V}(t, k)$  that is determined by the proposed procedure is equal to  $\tilde{V}(t, k)$ . Further the  $x_n$ 's are monotonically decreasing in  $n$ . □

## REFERENCES

- [1] AGRAWAL, V. and KAMBIL, A., “Dynamic pricing strategies in electronic commerce.” Working Paper, New York University and Accenture, 2000.
- [2] BAKER, K. R., “Sequencing rules and due-date assignments in a job shop,” *Management Science*, vol. 30, no. 9, pp. 1093–1104, 1984.
- [3] BAKER, K. and BERTRAND, J., “A dynamic priority rule for scheduling against due-dates,” *Journal of Operations Management*, vol. 3, no. 1, pp. 37–42, 1982.
- [4] BAKOS, Y. and BRYNJOLFSSON, E., “Bundling and competition on the internet,” *Marketing Science*, vol. 19, no. 1, p. 63, 2000.
- [5] BALAKRISHNAN, N., SRIDHARAN, S. V., and PATTERSON, J. W., “Rationing capacity between two product classes,” *Decision Sciences*, vol. 27, no. 2, pp. 185–214, 1996.
- [6] BELOBABA, P. P., “Airline yield management - an overview of seat inventory control,” *Transportation Science*, vol. 21, no. 2, pp. 63–73, 1987.
- [7] BELOBABA, P. P., “Application of a probabilistic decision-model to airline seat inventory control,” *Operations Research*, vol. 37, no. 2, pp. 183–197, 1989.
- [8] BERTSIMAS, D. and POPESCU, I., “Revenue management in a dynamic network environment,” *Transportation Science*, vol. 37, no. 3, pp. 257–277, 2003.
- [9] BILLESBACH, T., HARRISON, A., and CROOM-MORGAN, S., “Supplier performance measures and practices in JIT companies in the U.S. and the U.K.,” *International Journal of Purchasing and Materials Management*, vol. 27, pp. 24–29, 1991.
- [10] BOYACI, T. and RAY, S., “Product differentiation and capacity cost interaction in time and price sensitive markets,” *Manufacturing & Service Operations Management*, vol. 5, no. 1, pp. 18–36, 2003.
- [11] CARBONE, J., “Some independents try to reinvent themselves,” *Purchasing*, vol. June 18, 1998.
- [12] CHAN, L. M. A., SIMCHI-LEVI, D., and SWANN, J. L., “Pricing, production, and inventory policies for manufacturing with stochastic demand and discretionary sales,” *Manufacturing and Service Operations Management*, vol. 8, no. 2, pp. 149–168, 2006.
- [13] CHARNSIRISAKSKUL, K., GRIFFIN, P., and KESKINOCAK, P., “Order selection and scheduling with lead-time flexibility,” *IIE Transactions*, vol. 36, pp. 697–707, 2004.
- [14] CHARNSIRISAKSKUL, K., GRIFFIN, P., and KESKINOCAK, P., “Pricing and scheduling decisions with lead-time flexibility,” *European Journal of Operational Research*, vol. 171, no. 1, pp. 153–169, 2006.



- [15] CHATTERJEE, S., SLOTNICK, S., and SOBEL, M., “Delivery guarantees and the interdependence of marketing and operations,” *Production and Operations Management*, vol. 11, no. 3, pp. 393–409, 2002.
- [16] COURTY, P., “An economic guide to ticket pricing in the entertainment industry,” *Louvain Economic Review*, vol. 66, no. 1, pp. 167–192, 2000.
- [17] DEKKER, R., HILL, R. M., KLEIJN, M. J., and TEUNTER, R. H., “On the (s-1, s) lost sales inventory model with priority demand classes,” *Naval Research Logistics*, vol. 49, no. 6, pp. 593–610, 2002.
- [18] DELLAERT, N., “Due-date setting and production control,” *International Journal of Production Economics*, vol. 23, pp. 59–67, 1991.
- [19] DENECKER, R. and PECK, J., “Competition over price and service rate when demand is stochastic: a strategic analysis,” *RAND Journal of Economics*, vol. 26, no. 1, pp. 148–162, 1995.
- [20] DEPAOLI, L., “Executive Vice President & Chief Marketing Officer, Atlanta Spirit, LLC,” 2006.
- [21] DESHPANDE, V., COHEN, M. A., and DONOHUE, K., “A threshold inventory rationing policy for service-differentiated demand classes,” *Management Science*, vol. 49, no. 6, pp. 683–703, 2003.
- [22] DRAKE, M., DURAN, S., GRIFFIN, P., and SWANN, J., “A static model of revenue management for sports and entertainment tickets.” Working Paper, Georgia Institute of Technology, 2006.
- [23] DUENYAS, I., “Single facility due date setting with multiple customer classes,” *Management Science*, vol. 41, no. 4, pp. 608–619, 1995.
- [24] DUENYAS, I. and HOPP, W., “Quoting customer lead times,” *Management Science*, vol. 41, no. 1, pp. 43–57, 1995.
- [25] ENNS, S., “Lead time selection and the behaviour of work flow in job shops,” *European Journal of Operational Research*, vol. 109, pp. 122–136, 1998.
- [26] FARLEY, T., “Groups need revenue management too,” *Journal of Revenue & Pricing Management*, vol. 2, no. 2, p. 153, 2003.
- [27] FENG, Y., *Continuous-time Models in Perishable Asset Revenue Management*. PhD thesis, Columbia University, 1994.
- [28] FENG, Y. and GALLEGRO, C., “Perishable asset revenue management with markovian time dependent demand intensities,” *Management Science*, vol. 46, no. 7, pp. 941–956, 2000.
- [29] FENG, Y. and GALLEGRO, G., “Optimal starting times for end-of-season sales and optimal stopping-times for promotional fares,” *Management Science*, vol. 41, no. 8, pp. 1371–1391, 1995.
- [30] FENG, Y. and XIAO, B., “Maximizing revenues of perishable assets with a risk factor,” *Operations Research*, vol. 47, no. 2, pp. 337–341, 1999.

- [31] FLEISCH, E. and POWELL, S. G., “The value of information integration in meeting delivery dates,” *Journal of Organizational Computing and Electronic Commerce*, vol. 11, no. 1, pp. 15–30, 2001.
- [32] FRANK, K. C., ZHANG, R. Q., and DUENYAS, I., “Optimal policies for inventory systems with priority demand classes,” *Operations Research*, vol. 51, no. 6, pp. 993–1002, 2003.
- [33] FRY, T., PHILIPOOM, P., and MARKLAND, R., “Due date assignment in a multistage job shop,” *IIE Transactions*, pp. 153–161, June 1989.
- [34] GALLEGO, G. and PHILLIPS, R., “Revenue management of flexible products,” *Manufacturing & Service Operations Management*, vol. 6, no. 4, pp. 321–337, 2004.
- [35] GALLEGO, G. and VAN RYZIN, G., “Optimal dynamic pricing of inventories with stochastic demand over finite horizons,” *Management Science*, vol. 40, no. 8, pp. 999–1020, 1994.
- [36] GANS, N., “Customer loyalty and supplier quality competition,” *Management Science*, vol. 48, no. 2, pp. 207–221, 2002.
- [37] GERAGHTY, M. K. and JOHNSON, E., “Revenue management saves National Car Rental,” *Interfaces*, vol. 27, no. 1, pp. 107–127, 1997.
- [38] GUPTA, D. and WANG, L., “Manufacturing capacity revenue management,” *Forthcoming in Operations Research*, 2006.
- [39] HA, A. Y., “Inventory rationing in a make-to-stock production system with several demand classes and lost sales,” *Management Science*, vol. 43, no. 8, pp. 1093–1103, 1997.
- [40] HA, A. Y., “Stock rationing in an m/e-k/1 make-to-stock queue,” *Management Science*, vol. 46, no. 1, pp. 77–87, 2000.
- [41] HALL, J. and PORTEUS, E., “Customer service competition in capacitated systems,” *Manufacturing & Service Operations Management*, vol. 2, no. 2, p. 144, 2000.
- [42] HAUSMAN, W. H., MONTGOMERY, D. B., and ROTH, A. V., “Why should marketing and manufacturing work together? some exploratory empirical results,” *Journal of Operations Management*, vol. 20, no. 3, pp. 241–257, 2002.
- [43] HOF, R. D., SAGER, I., and HIMELSTEIN, L., “The sad saga of silicon graphics,” *Business Week Online, Cover Story*, 1997.
- [44] HOPP, W. and ROOF STURGIS, M., “Quoting manufacturing due dates subject to a service level constraint,” *IIE Transactions*, vol. 32, pp. 771–784, 2000.
- [45] KARAESMEN, I. and VAN RYZIN, G., “Overbooking with substitutable inventory classes,” *Operations Research*, vol. 52, no. 1, pp. 83–104, 2004.
- [46] KARATZAS, I. and SHREVE, S., *Brownian Motion and Stochastic Calculus*. NY: Springer-Verlag, 1988.

- [47] KATIRCIOGLU, K. and ATKINS, D., “Managing inventory with multiple customer classes requiring different levels of service.” Working Paper, IBM, 1996.
- [48] KESKINOCAK, P., RAVI, R., and TAYUR, S., “Scheduling and reliable lead time quotation for orders with availability intervals and lead time sensitive revenues,” *Management Science*, vol. 47, no. 2, pp. 264–279, 2001.
- [49] KESKINOCAK, P. and TAYUR, S., *Due-Date Management Policies*. in Handbook of Quantitative Supply Chain Analysis: Modeling in the E-Business Era, D. Simchi-Levi, S.D. Wu, and Z.M. Shen (editors), Kluwer Academic Publishers, 2004.
- [50] KRANTZ, M., “Victors, vanquished in year’s rough waters,” *USA Today*, 2000.
- [51] LESLIE, P., “Price discrimination in Broadway theater,” *RAND Journal of Economics*, vol. 35, no. 3, pp. 520–541, 2004.
- [52] LIEBERMAN, W., “Implementing yield management,” in *ORSA/TIMS National Meeting*, (San Francisco, California), 1992.
- [53] LIEBESKIND, J. and RUMELT, R., “Markets for experience goods with performance uncertainty,” *RAND Journal of Economics*, vol. 20, no. 4, pp. 601–621, 1989.
- [54] LITTLEWOOD, K., “Forecasting and control of passenger bookings,” *AGIFORS Symposium Proceedings*, pp. 95–117, 1972.
- [55] LIU, T. and SIMCHI-LEVI, D., “Delayed production strategies with backlogged and discretionary sales.” Working Paper, MIT, 2003.
- [56] MCAFEE, R. P., MCMILLAN, J., and WHINSTON, M. D., “Multiproduct monopoly, commodity bundling, and correlation of values,” *Quarterly Journal of Economics*, vol. 104, no. 2, pp. 371–383, 1989.
- [57] MELCHORS, P., DEKKER, R., and KLEIJN, M. J., “Inventory rationing in an  $(s, q)$  inventory model with lost sales and two demand classes,” *Journal of the Operational Research Society*, vol. 51, no. 1, pp. 111–122, 2000.
- [58] MENDELSON, H. and WHANG, S., “Optimal incentive-compatible priority pricing for the  $m/m/1$  queue,” *Operations Research*, vol. 38, no. 5, p. 870, 1990.
- [59] MERCHANT REVIEWS [http://shopping.yahoo.com/merchrating/user\\_rv.html](http://shopping.yahoo.com/merchrating/user_rv.html), last accessed on July 27, 2006.
- [60] MIYAZAKI, S., “Combined scheduling system for reducing job tardiness in a job shop,” *International Journal of Production Research*, vol. 19, no. 2, pp. 201–211, 1981.
- [61] MOON, I. and KANG, S., “Rationing policies for some inventory systems,” *Journal of the Operational Research Society*, vol. 49, no. 5, pp. 509–518, 1998.
- [62] PALAKA, K., ERLEBACHER, S., and KROPP, D. H., “Lead-time setting, capacity utilization, and pricing decisions under lead-time dependent demand,” *IIE Transactions*, vol. 30, no. 2, pp. 151–163, 1998.
- [63] PEKGUN, P., GRIFFIN, P., and KESKINOCAK, P., “Coordination of marketing and production for price and leadtime decisions,” *Forthcoming in IIE Transactions*, 2006.

- [64] RAY, S. and JEWKES, E., “Customer lead time management when both demand and price are lead time sensitive,” *European Journal of Operational Research*, 2003.
- [65] ROGERS, L. and WILLIAMS, D., *Diffusions, Markov Processes and Martingales, Volume 2: Itô Calculus*. NY: John Wiley & Sons, 1987.
- [66] ROSEN, S. and ROSENFELD, A., “Ticket pricing,” *Journal of Law & Economics*, vol. 40, no. 2, pp. 351–376, 1997.
- [67] SALINGER, M. A., “A graphical analysis of bundling,” *Journal of Business*, vol. 68, no. 1, pp. 85–98, 1995.
- [68] SCARF, H. E., “Optimal inventory policies when sales are discretionary.” Cowles Foundation Discussion Paper No. 1270, 2000.
- [69] SLOTNICK, S. A. and SOBEL, M. J., “Manufacturing lead-time rules: Customer retention versus tardiness costs,” *European Journal of Operational Research*, vol. 163, no. 3, pp. 825–856, 2005.
- [70] SMITH, B. C., LEIMKUEHLER, J. F., and DARROW, R. M., “Yield management at American Airlines,” *Interfaces*, vol. 22, no. 1, pp. 8–31, 1992.
- [71] SO, K., “Price and time competition for service delivery,” *Manufacturing & Service Operations Management*, vol. 2, no. 4, pp. 392–409, 2000.
- [72] SO, K. and SONG, J.-S., “Price, delivery time guarantees and capacity selection,” *European Journal of Operational Research*, vol. 111, pp. 28–49, 1998.
- [73] SOBEL, M. J. and ZHANG, R. Q., “Inventory policies for systems with stochastic and deterministic demand,” *Operations Research*, vol. 49, no. 1, pp. 157–162, 2001.
- [74] SPEARMAN, M. and ZHANG, R., “Optimal lead time policies,” *Management Science*, vol. 45, no. 2, pp. 290–295, 1999.
- [75] SPECIAL REPORT <http://www.icaew.co.uk/index.cfm?route=136544>, last accessed on July 27, 2006.
- [76] STIDHAM, S., “Socially and individually optimal-control of arrivals to a GI/M/1 queue,” *Management Science*, vol. 24, no. 15, pp. 1598–1610, 1978.
- [77] STIDHAM, S., “Pricing and capacity decisions for a service facility: Stability and multiple local optima,” *Management Science*, vol. 38, no. 8, pp. 1121–1139, 1992.
- [78] SUPPLIER SCORECARD [http://www.samtec.com/standard\\_products/quality\\_information](http://www.samtec.com/standard_products/quality_information), last accessed on July 27, 2006.
- [79] TALLURI, K. and VAN RYZIN, G., “An analysis of bid-price controls for network revenue management,” *Management Science*, vol. 44, no. 11, p. 1577, 1998.
- [80] TOPKIS, D. M., “Optimal ordering and rationing policies in a nonstationary dynamic inventory model with  $n$  demand classes,” *Management Science*, vol. 15, pp. 160–178, 1968.

- [81] TOPKIS, D. M., *Supermodularity and Complementarity*. Princeton university Press, 1998.
- [82] VAN MIEGHEM, J., “Price and service discrimination in queuing systems: Incentive compatibility of gcmu scheduling,” *Management Science*, vol. 46, no. 9, p. 1249, 2000.
- [83] VEINOTT, A. F., “Optimal policy in a dynamic single product nonstationary inventory model with several demand classes,” *Operations Research*, vol. 13, pp. 761–778, 1965.
- [84] VENKATESH, R. and KAMAKURA, W., “Optimal bundling and pricing under a monopoly: Contrasting complements and substitutes from independently valued products,” *Journal of Business*, vol. 76, no. 2, pp. 211–231, 2003.
- [85] VENKATESH, R. and MAHAJAN, V., “A probabilistic approach to pricing a bundle of products or services,” *Journal of Marketing Research (JMR)*, vol. 30, no. 4, p. 494, 1993.
- [86] WEEKS, J., “A simulation study of predictable due-dates,” *Management Science*, vol. 25, no. 4, pp. 363–373, 1979.
- [87] WEIN, L., “Due-date setting and priority sequencing in a multi class M/G/1 queue,” *Management Science*, vol. 37, no. 7, pp. 834–850, 1991.
- [88] WILLIAMSON, E., *Airline network seat control*. PhD thesis, MIT, 1992.
- [89] YUEN, B. B., “Group revenue management: Redefining the business process — part I,” *Journal of Revenue & Pricing Management*, vol. 1, no. 3, p. 267, 2002.

## VITA

Serhan Duran was born in Türkoğlu, Turkey on October 2, 1980. He received his B.S. in Industrial Engineering from Middle East Technical University in 2002. He started pursuing his Ph.D. in the School of Industrial and Systems Engineering in August 2002 at Georgia Institute of Technology. He received an M.S. in Industrial Engineering and an M.S. in Operations Research from the Georgia Institute of Technology in 2004 and 2005, respectively. His research interests include using stochastic models for the improvement of manufacturing, service and distribution operations to increase flexibility. He will be joining the faculty of the Industrial Engineering Department at Middle East Technical University, in Fall 2007.