

Essays on Retail Management with Emerging Practice and Customer Behavior

by

Lai Wei

A dissertation submitted in partial fulfillment
of the requirements for the degree of
Doctor of Philosophy
(Business Administration)
in The University of Michigan
2020

Doctoral Committee:

Associate Professor Stefanus Jasin, Co-Chair
Professor Roman Kapuscinski, Co-Chair
Assistant Professor Heng Liu
Associate Professor Brian Wu

Lai Wei

laiwi@umich.edu

ORCID iD: [0000-0001-5451-5059](https://orcid.org/0000-0001-5451-5059)

© Lai Wei 2020

TABLE OF CONTENTS

LIST OF FIGURES	iv
LIST OF TABLES	v
LIST OF APPENDICES	vi
ABSTRACT	vii
CHAPTER	
I. Shipping Consolidation across Two Warehouses with Delivery Deadline and Expedited Options for E-commerce and Omni- channel Retailers	
1.1 Introduction	1
1.2 Brief Literature Review	8
1.3 Model Specification	11
1.4 Single Warehouse with Only Fixed Cost	15
1.4.1 The Optimal Policy: Its Structure and Properties	16
1.5 Two Warehouses with Only Fixed Cost	19
1.5.1 The Optimal Policy: Its Structure and Properties	20
1.5.2 Simple Heuristics	23
1.6 One Warehouse with Fixed and Variable Costs	28
1.6.1 Dynamic Programming (DP) Formulation	28
1.6.2 A Simple Heuristic	31
1.7 Two Warehouses with Fixed and Variable Costs	34
1.7.1 Base Heuristics	35
1.7.2 Performance of Base Heuristics: Warehouse-Based and Order-Based	35
1.7.3 Warehouse-based ⁺⁺ Heuristic	39
1.8 Conclusion	44

II. On a Deterministic Approximation of Inventory Systems with Sequential Probabilistic Service Level Constraints	46
2.1 Introduction	46
2.1.1 Outline of Paper	51
2.2 Model Description	51
2.3 Backorder Inventory System	57
2.3.1 Proposed Heuristic Control and Its Performance . . .	58
2.3.2 Proof of Theorem II.5	60
2.3.3 Numerical Results	64
2.4 Lost-Sales Inventory System	67
2.4.1 Proposed Heuristic Control and Its Performance . . .	67
2.4.2 Proof of Theorem II.9	69
2.5 Conclusion	78
III. Snob and Follower Effects in Luxury Retailing	80
3.1 Introduction	80
3.2 Literature Review	86
3.3 Model	89
3.4 Product-line Strategy	91
3.5 Selling Strategy in Competitive Market	98
3.5.1 No Externalities	99
3.5.2 Snob and Follower Effects	104
3.6 Bundling Strategy	117
3.7 Conclusions	125
APPENDICES	127
BIBLIOGRAPHY	178

LIST OF FIGURES

Figure

1.1	Boundaries for type A orders when $z_B = 5$	21
1.2	Boundary τ_B^1 for type B orders when $z_A = 11$	21
1.3	Asymmetric case 1 when $z_A = 8$	23
1.4	Asymmetric case 2 when $z_C = 9$	23
1.5	Average percentage gaps across different d , α_A , and α_C	41
1.6	Average percentage gaps across different d and α_C , with $\alpha_B = 0.9$	42
3.1	Optimal Policies without Externalities	123
3.2	Optimal Policies with Externalities	124

LIST OF TABLES

Table

1.1	Percentage Gaps to Optimal Cost	33
1.2	Percentage Gap to Optimal	37
1.3	Percentage Gap to Order-Based Heuristic	37
1.4	Percentage Gap to Optimal	40
1.5	Percentage Gap to Order-Based Heuristic	40
2.1	Percentage Performance Gap (in %) with Poisson Distributed Demand	65
2.2	Percentage Performance Gap (in %) with Normal Distributed Demand	66
3.1	6 Cases of Demand Functions	109
3.2	Equilibriums in Competitive Market	114

LIST OF APPENDICES

Appendix

A.	Parameters in Simulation Experiments of Chapter I	128
B.	Proof of Chapter I	129
C.	Proof of Chapter II	147
D.	Proof of Chapter III	163

ABSTRACT

This dissertation is motivated by observations of emerging retail practice and the corresponding customer requirements and behavior. Specifically, it includes the following topics: (1) omni-channel and e-commerce logistics, (2) inventory management, and (3) product strategy and pricing analytics.

In the first chapter of the dissertation, we study the shipment consolidation policy facing an increasing frequency of orders per customer. Shipment consolidation (i.e., shipping multiple orders together instead of shipping them separately) is commonly used to decrease total shipping costs. However, when the delivery of some orders is delayed so they can be consolidated with future orders, a more expensive expedited shipment may be needed to meet shorter deadlines. In this paper, we study the optimal consolidation policy focusing on the trade-off between economies of scale due to combining orders and expedited shipping costs, in the setting of two warehouses. Our work is motivated by the application of fulfillment consolidation in e-commerce and omni-channel retail, especially with the rise of so-called on-demand logistics services. Sellers have the flexibility to take advantage of consolidation by deciding when to ship the orders and from which warehouse to fulfill the orders, as long as the orders' deadlines are met. The optimal policies and their structures are characterized. Using the insights of these structural properties, we propose two easily implementable heuristics that perform within 1-2% of the optimal solution and outperform other benchmark consolidation methods in numerical tests.

In the second chapter of the dissertation, we study the inventory decision when there are explicit high service-level requirements. We consider a stochastic inventory model (under both backorder and lost-sales) with non-stationary demands, positive lead times, and sequential probabilistic service level constraints. This is a notoriously difficult problem to solve and, to date, not much progress has been made in understanding the structure of its optimal control, especially for the lost-sales inventory system. In this paper, we propose a simple order-up-to control, whose parameters can be calculated using the optimal solution of a deterministic approximation of the backorder inventory system, and show that it is asymptotically optimal for both the backorder and lost-sales systems in the regime of high service level requirement. This result contributes to the growing body of inventory literature that show the near-optimality of simple heuristic controls. Moreover, it also gives credence to the use of deterministic approximation for solving complex inventory problems in practice, at least for applications where the targeted service level is sufficiently high.

In the third chapter of the dissertation, we study product strategy and pricing analytics, in settings where customers have both positive and negative product network externalities. One unique feature of luxury products is the coexistence of two opposite externalities: snob customers experience negative externalities with product sales while follower customers experience positive externalities. Motivated by several interesting and (perhaps) counter-intuitive practices in the luxury industry, we study the effect of these two opposite externalities with respect to the selling strategies from three perspectives: 1) the product-line strategy in a monopoly setting, 2) the pricing strategy in a competition setting, and 3) the product bundling strategy. We find that these two opposite externalities generally work in the same direction, although through different mechanisms.

CHAPTER I

Shipping Consolidation across Two Warehouses with Delivery Deadline and Expedited Options for E-commerce and Omni-channel Retailers

1.1 Introduction

The total costs of logistics usually account for 9% to 14% of sales of a company, depending on the industry sector and, at an aggregate level, represent 7.9% of the US GDP in 2015. Among them, shipping costs alone comprise more than 60% (27th State of Logistics Report). For a typical online retailer, Amazon.com, the shipping costs account for as high as 11.89% of the net sales (Amazon 2016 Annual Report). Thus, it is no surprise that “effectively managing shipping costs directly affects . . . business’ bottom line” (Fell, 2011). One commonly used strategy to save on shipping costs is to consolidate multiple small shipments into a large one (Cetinkaya, 2005). Although the value of consolidation is already well recognized in the supply chain, its saving potential in the Business-to-Customer (B2C) setting, especially for e-commerce and omni-channel retailers, has not been fully exploited yet and multiple service providers continue experimentations in this area. Retailers often have significant opportunities to take advantage of various forms of flexibility when satisfying customers’ orders,

especially choosing from *which* warehouse to fulfill the orders and *when* to ship the orders. In this paper, we focus on the latter, taking advantage of a time window between the time the retailer receives an order and the time by which the order needs to be delivered (Lee et al., 2001; Xu et al., 2009). This window provides an opportunity to combine existing orders with new incoming orders (either from the same customer or multiple customers located in nearby regions) and ship them together. However, since orders must be delivered by their guaranteed due dates (or deadlines), delaying the shipment of some orders may increase total shipping costs due to the need to use expedited shipping. Many logistics firms such as Expedited Logistics and Freight Services¹, ASAP Expedited Logistics², etc., increase shipping rates for faster shipping modes. Major carriers such as UPS³, FedEx⁴, and USPS⁵ also offer shorter lead time deliveries (3 Day Select, Overnight Delivery, etc.) for higher fees.

Our work is motivated by the potential benefit of consolidation in e-commerce and omni-channel settings. For e-commerce, the frequency of orders per customer has continued to increase while the size of each order has become smaller, partly due to the popularity of free-shipping service offered by retailers such as Neimanmarcus.com and others (Lewis, 2006; Gil, 2014). In this domain, consolidation can be, and is, implemented for individual customers. Indeed, one of the co-founders of one of the largest e-commerce retailers in China recently told us that the proportion of its customers who place two orders within an hour and one day is 5% and 20%, respectively. In addition, the proportion of customers who place three or more orders within one day and one week is 1% and 10%, respectively. Given that many orders are shipped for free, this high customer ordering frequency has intensified pressure to take advantage

¹www.elfsfreight.com/services-domestic.php

²www.asapexpediting.net

³www.ups.com/us/en/shipping/zones-and-rates/48-contiguous-states.page

⁴images.fedex.com/us/services/pdf/FedEx_StandardListRates_2018.pdf

⁵pe.usps.com/text/dmm300/Notice123.htm

of such situations. The consolidation of multiple orders placed by the same customer can bring significant savings, due to the “fixed-cost” portion of the shipping costs—we observe that the cost structures of the main third-party logistics firms in the US and in China (UPS, USPS, FedEx, SF Express⁶, Yunda Express⁷) are all in the form of fixed cost + variable cost × volume, where the fixed cost is usually more than five times of the unit variable cost. Thus, consolidating multiple orders avoids incurring the fixed cost multiple times. Indeed, the firm we talked to expressed concerns about missed opportunities in reducing costs and considered testing order consolidation: “Under our current practice, multiple orders from a customer are dispatched individually. Holding orders for a longer time, especially for customers who have a record of frequently impulsive purchases, reduces the number of deliveries, thus saves costs. Such cost reduction may well compensate more expensive expedited delivery service. Consolidation would be strongly considered as a future step in our operations improvements.” (Anonymous, Personal Communication, 2017)

Aside from the e-commerce setting, our work is also directly applicable in omnichannel setting, for consolidating orders from multiple customers. While having a 5-10 days delivery window was typical a decade ago, today’s standard is a lot more aggressive and many retailers now offer either a one-day or several-hour delivery guarantee from their stores (Hausmann, 2014). To meet the same-day deadline as well as to control for the total shipping costs, the retailers increasingly rely on the third-party-logistics firms and fulfill orders from the nearest omnichannel stores.⁸ Given that the third-party-logistics firms charge a fixed cost for any delivery plus additional variable costs for each extra stop (usually much smaller than the fixed cost, as we explained below), it is

⁶www.sf-express.com/cn/en/dynamic-function/price

⁷www.yundaex.com/cn

⁸Many third-party logistics companies (e.g., Pulse-commerce and Onestock), who arrange the ship-from-store for omnichannel retailers, also list shipping from nearest store as their policy. More details can be found at www.pulse-commerce.com/omnichannel-retail-ship-from-store-slideshare and onestock-retail.com/en/ship-from-store

beneficial for the retailer to consolidate orders from customers living in nearby regions and dispatch one delivery with multiple drop-offs. The following example illustrates the potential savings: Consider typical order deliveries from downtown areas to midtown areas in NYC using Breakaway Courier, a NYC-based on-demand logistic firm. The firm provides several expedited options: 90-min, 60-min, 40-min, and 20-min deliveries, with the corresponding multipliers of $1\times$, $1.5\times$, $2\times$, $3\times$ times regular price (both fixed and variable costs), respectively. For the 90-min delivery option, the fixed cost (per delivery) is \$24.95 and the regular variable cost (per additional stop) is \$7.25. Thus, sending three separate orders with the 90-min service incurs a total costs of $3\times\$24.95 = \74.85 whereas sending them together incurs a total costs of $\$24.95 + 2\times\$7.25 = \$39.45$ (a \$35.4 saving). If we send all three orders together using the 60-min service, we incur a total costs of $1.5\times(\$24.95 + 2\times\$7.25) = \$59.175$ (a \$15.675 saving, in case the second order arrives so late that a faster service needs to be used). Obviously, an even larger total savings can be achieved when the number of customers who place such request (i.e., the number of orders that can be consolidated) is higher.

In this paper, we study the optimal shipping and consolidation policy by taking into account both the delivery deadline and the expedited shipping options. We assume that the retailer operates up to two warehouses/stores (i.e., the primary warehouses, or local stores, for some regions). This reflects the current practice of many online and omni-channel retailers as they typically ship orders from two nearest warehouse for one customer (Acimovic, 2015).⁹ To take advantage of consolidation, the retailer may either consolidate orders from the same customer or from multiple customers living in nearby regions. Each shipment incurs both fixed and variable costs, whose

⁹In the omni-channel situations we consider in the paper, where the deadline is within hours or one day, there is not much room for the retailers to consider stores other than the nearby ones. Third-party service-providing companies, who arrange the ship-from-store for omni-channel retailers, also list shipping from nearest store as their policy, e.g. Pulsecommerce (www.pulse-commerce.com/omnichannel-retail-ship-from-store-slideshare) and Onestock (onestock-retail.com/en/ship-from-store).

values depend on the delivery speed. We consider a finite-horizon problem where, in each period, the retailer needs to make three joint decisions: (1) *Which orders should be shipped?* (2) *From which warehouse/store should the orders be shipped?* (3) *How should the shipment be split into multiple packages/delivery tasks?* The first decision is illustrated in the two examples of e-commerce and omni-channel retailers above. We illustrate the second decision using the cost structure of Breakaway Courier. Consider a simple setting where a retailer operates two stores/warehouses located in downtown NYC and fulfills orders from customers located in the midtown. The retailer has three pending orders, the first order is due in 90 minutes and can be fulfilled from store/warehouse 1, the second order is due in 60 minutes and can be fulfilled from both stores/warehouses, and the third order is due in 60 minutes and can be fulfilled from store/warehouse 2. Consider the case where we know that a new order will arrive in 20 minutes (no other orders will arrive on the day) and that this order can only be fulfilled from store/warehouse 2. In this case, the first order should be shipped immediately, as there is no consolidation opportunity. However, it is not obvious whether the second order should be consolidated with the first one. One alternative is to send both the first and second orders in the current period, which incurs a shipping cost of $1.5 \times (\$24.95 + \$7.25) = \$48.3$, and send the third order together with the new order after 20 mins, which incurs a shipping cost of $2 \times (\$24.95 + \$7.25) = \$64.4$. The total shipping costs is \$112.7. Another alternative is to send only the first order in the current period, which incurs a shipping cost of \$24.95, and send all the remaining orders together after 20 mins, which incurs a shipping cost of $2 \times (\$24.95 + 2 \times \$7.25) = \$78.9$. The total shipping costs of this alternative is \$103.85. This example illustrates the complexity of consolidation decision even in the setting with only two warehouses, as it is not always economical to consolidate existing orders.

The third decision that the retailers need to make is package splitting. It is not

necessarily optimal to ship all orders in one package. Splitting the very urgent orders from the less urgent ones allows us to avoid incurring high variable costs (due to the need to use an urgent shipping mode) for the less urgent ones as the less urgent orders can be shipped using slower and cheaper shipping mode. Taking into account future hypothetical arrivals makes the decision even less obvious. This decision is discussed in more details in Section 1.3 and 1.6.

Given the complexity of the problem (see Section 1.3), after providing the general formulation of the problem, we proceed by analyzing three simplified cases: one-warehouse setting with only fixed cost, two-warehouse setting with only fixed cost, and one-warehouse setting with both fixed and variable costs. The insights from these cases are then used to construct easy-to-implement heuristics for two-warehouse setting with both fixed and variable costs. Our findings in this paper can be summarized as follows:

1. For the one-warehouse setting with only fixed cost, we show that the optimal policy can be characterized by a sequence of time-dependent thresholds—it is optimal to ship all pending orders in period t if the *slack time* (remaining time until the deadline) of the most urgent order is smaller than or equal to a threshold τ_t . This result is intuitive: Retailers can take advantage of consolidation to the point where the increase in total costs becomes so high that it exceeds the potential benefit of consolidation.
2. For the two-warehouse setting with only fixed cost, we show that, in general, the optimal policy is no longer easy to characterize. For the special case where the two warehouses are symmetric, the optimal policy can be characterized by six non-linear boundaries in a three-dimensional space. Motivated by the simplicity of threshold policy in the one-warehouse setting, we propose two heuristics that replace these six boundaries with constant thresholds: *warehouse-based* heuristic and *order-based* heuristic. Under the warehouse-based heuristic, once the

threshold for a warehouse is crossed, all orders that are “shippable” from that warehouse are shipped; under the order-based heuristic, once the threshold for an order *type* is crossed, all orders of that type are shipped, together with some other orders that can be consolidated. Our numerical experiments, based on the typical range of customer ordering frequency and on the structure of UPS shipping rates, show that the performances of these heuristics across symmetric and asymmetric problem instances are within 2% of the optimal policy in most cases.

3. For the one-warehouse setting with both fixed and variable costs, we show that adding variable costs into the model makes the analysis significantly more complicated compared to that of the setting with only fixed cost—more details about this complexity can be found in Section 1.3. That said, for the special case where all orders have guaranteed delivery within at most three periods, we are able to show that the optimal policy can be characterized by thresholds that are a function of order volume.
4. For the two-warehouse setting with both fixed and variable costs, we do not attempt to analyze the structure of the optimal policy due to the already complex structure of the optimal policy in the setting with only fixed cost. Instead, given the good performance of constant-threshold heuristics in the previous settings, we propose modified heuristics and show, using numerical study, that their average performances are within 0.29%-1.83% of the optimal policy. While, in most of the cases, the heuristics perform well, we also identify two conditions where the average optimality gap is twice as large as that in the other cases. Such poor performance happens when the mismatch between the pre-determined allocation, when calculating the thresholds, and the state-dependent allocation, when utilizing the thresholds, is exaggerated. However, the parameters that meet these

two conditions are extreme and correspond to a setting that rarely occurs in reality. Thus, despite the complexity of the actual optimal policy, simple heuristics perform well in most cases.

To conclude, by considering setting that capture key elements of real life situations, we have illustrated that consolidation is an effective way to improve the standard outbound shipping policy, with significant cost reductions due to anticipating and properly planning consolidation across time. Such policies can be implemented by applying the threshold-form heuristics we propose.

The remainder of the paper is organized as follows: Section 2 provides a brief literature review. In Section 1.3, we provide a general formulation of the problem. In Sections 1.4, 1.5, and 1.6, we study the three simplified cases. In Section 1.7, we study the two-warehouse setting with both fixed and variable costs and numerically test the proposed heuristics. In Section 1.8, we conclude the paper. All proofs can be found in the supplemental file.

1.2 Brief Literature Review

Two streams of literature are most closely related to our work: shipment consolidation, which studies how to combine several orders, and order fulfillment, which studies from which warehouse to fulfill the orders. The potential cost savings due to shipment consolidation have been extensively studied in the logistics literature (Daganzo, 1988; Pooley and Stenger 1992; Popken, 1994). The main trade-off considered in this literature is between the constant fixed cost of shipping and the inventory holding cost. Three types of consolidation policies are usually considered: (1) time-based, which sets a pre-determined interval within which orders are accumulated and one shipment is dispatched at the end of the interval; (2) quantity-based, which dis-

patches one shipment after a pre-determined quantity of orders is accumulated; and (3) hybrid, or time-and-quantity, consolidation, which releases a shipment either after a pre-determined quantity is achieved or at the end of a pre-determined time interval. All these policies are heuristics—most existing literature either focuses on evaluating and comparing the performance of these policies (Cooper, 1984; Burns et al., 1985; Campbell, 1990; Higginson and Bookbinder, 1994) or on calculating their optimal parameters with or without integrating inventory decisions (Gupta and Bagchi, 1987; Axsater, 2001; Cetinkaya et al., 2000, 2008; Popken 1994). Our work differs from the previous consolidation literature in two ways: (1) We are the first to study shipping consolidation in the context of e-commerce and omni-channel retailing. Unlike in the setting considered in the existing literature, where the main trade-off is between the fixed cost and the holding cost, in our setting, orders are held for at most a few hours to a few days, which means that the extra holding cost due to waiting is negligible and the main trade-off is between the fixed cost and the additional speed-up cost, due to expedited shipment.¹⁰ Although expedited shipment has been considered in the inventory literature (Zhou and Chao, 2010; Caggiano, etc., 2006; Huggins and Olsen, 2003; Hoadley and Heyman, 1977) and supply chain risk management (Qi and Lee, 2015), it has not been considered in the context of consolidation. (2) Unlike most works in the existing literature that focus primarily on heuristics, in this paper, we analyze the structure of the optimal policy (at least for some cases) and then propose some heuristics based on the structure of this optimal policy. This approach yields a new insight as our heuristics do not follow the same structure as in the commonly-used

¹⁰Note that our model and the holding-cost model are not mathematically equivalent. This is because that while the unit holding cost and fixed cost in consolidation vs. holding cost setting is constants, the unit variable cost and fixed shipping costs we considered in the paper depend on the orders' remaining time to the deadline. Thus, our problem has different structure and cannot be translated into the holding-cost problem and actually is much more complex.

time/quantity-based consolidation policy.¹¹

Another stream of relevant literature analyzes order fulfillment. Order fulfillment in the context of e-commerce retail has recently received a lot of attention. The main trade-off considered in this stream of literature is between shipping costs and future product availability, since not all warehouses stock the same products and there is typically inventory imbalance across different warehouses.¹² Xu et al. (2009) consider the delivery time window in designing a shipment policy, but their focus is on the fulfillment of individual orders instead of the consolidation of current order with future orders. Acimovic and Graves (2015) and Jasin and Sinha (2015) further analyze the model introduced in Xu et al. (2009) by designing tractable near-optimal shipment heuristics. They consider a very general model with multiple warehouses, multiple customer locations, and shipping costs that depend explicitly on the shipping distances. Lei et al. (2016) analyze the joint pricing and fulfillment problem. Since both pricing and fulfillment jointly affect the distribution of demand and supply across the system, they argue that these decisions must be considered jointly. None of these papers, however, considers consolidation of orders arriving at different times. While our work focuses on a new and economically relevant consolidation, we do not address the problem of inventory balancing across different warehouses in the current paper. Combining consolidation and future inventory availability in a single model is analytically intractable as even the time consolidation part itself is already not trivial (as we show in this paper). We consider these two as complementary approaches, as there are settings in which future inventory availability is more crucial than considering consolidation (e.g., when replenishments are infrequent) and vice versa (e.g., when replenishments are very

¹¹Specifically, although all heuristics mentioned above are based on some kind of thresholds, we are not aware of any work in the literature that uses slack time to define a threshold. The definition of slack time is given in Section 3.

¹²A retailer may ship an order from a warehouse with higher shipping cost, hoping that a future order (containing multiple items) can be shipped from the warehouse in which the inventory was reserved.

frequent).

1.3 Model Specification

In this section, we first discuss a general model that incorporates the key elements reflecting the economy of shipping consolidation problem investigated in this paper. As shown below, the general model is hard to analyze directly. This motivates us to analyze special cases and propose heuristics for the general model in the subsequent sections.

We consider a finite-horizon problem with T periods where time is indexed backward with period 1 being the last period and period T being the first period. Orders arrive from either a single customer or a single region.¹³ Each incoming order must be delivered no later than d periods after its arrival.¹⁴ We define *slack time*, s , as the remaining time until the deadline, e.g., $s = 1$ means that the order must reach customer in the next period. Reflecting the current practice, where retailers usually ship orders from at most two closest warehouses/stores, we assume that the retailer operates up to two warehouses, W1 and W2. Since not all products are available in all warehouses (Xu et al., 2009),¹⁵ incoming orders are classified into three types: type A can only be fulfilled from W1, type C can only be fulfilled from W2, and type B can be fulfilled from either W1 or W2.¹⁶ We denote the set of pending orders by a vector of their corresponding

¹³Since the decision for each customer/region is separate, the problem trivially extends to multiple customers/regions.

¹⁴In practice, retailers usually promise a default fixed length of deadline for orders, e.g., two-day delivery when shipping from warehouses for Amazon prime. There may be alternatives within the same retailer—Amazon customers may choose delayed delivery for a credit, e.g. \$1. We do not consider multiple-delivery options in our paper.

¹⁵This may be due to either temporary stockout or policies that dictate product allocation across warehouses. Our conversations with operations managers of Target Corp indicate that although some products are held in all warehouses, there are also products that are stored only in a subset of them.

¹⁶Note that this general model also covers some special cases. For example, the nested product allocation policy, where one warehouse carries a subset of products of the other warehouse, is a special case where types C products do not exist.

slack time (we use ∞ to denote the case of no pending orders). We allow multiple orders to arrive in one period and, for mathematical convenience, we assume that each order contains exactly one item, with one unit of volume; if an order contains multiple items, we simply treat them as separate orders. Suppose that there are currently n_A , n_B , and n_C pending orders of type A , B , and C , respectively. Then, the corresponding slack times are denoted by $\vec{s} = (\vec{s}_A, \vec{s}_B, \vec{s}_C)$, where $\vec{s}_A, \vec{s}_B, \vec{s}_C$ are each a list of actual slack times of orders of type $X \in \{A, B, C\}$, sequenced from the smallest slack time to the largest. That is, $\vec{s}_X = (s_{X,1}, s_{X,2}, \dots, s_{X,n_X})$ and $s_{X,1} \leq s_{X,2} \leq \dots \leq s_{X,n_X}$ ($X \in \{A, B, C\}$), where n_X is the number of orders of type X . For the shipping costs, we use $F_i(p)$ and $v_i(p)$ to denote the fixed cost and unit variable cost of delivering an order in p periods from warehouse i . We assume that both F_i and v_i are non-increasing and convex functions, which is consistent with data observed in practice.¹⁷ Moreover, $F_i(0) = \infty$, $F_i(\infty) = 0$ and $v_i(0) = \infty$, $v_i(\infty) = 0$, i.e., all delivery times are positive and all orders must be delivered on time.¹⁸ For each package/delivery task, the incurred fixed and unit variable costs are determined by the slack time of the most urgent order in the package/delivery task, while the total variable costs are also determined by the number of orders. For example, if l orders with slack times $(\theta_1, \theta_2, \dots, \theta_l)$, where $\theta_1 \leq \theta_2 \leq \dots \leq \theta_l$, are shipped in one package from warehouse i ($i \in \{1, 2\}$), then the total shipping costs is $F_i(\theta_1) + l \cdot v_i(\theta_1)$.

The sequence of events is as follows: At the beginning of period t , the retailer first observes the slack times of all pending orders $\vec{s} = (\vec{s}_A, \vec{s}_B, \vec{s}_C)$ and then makes the

¹⁷The cost functions we use in the paper are based on the official (non-discounted) UPS rates for different shipping modes. We looked for sufficiently good statistical fit in terms of both the functional form and the parameters. The data shows that both the fixed and variable shipping cost are non-decreasing and convex in the length of the delivery window. Translating these into slack times is natural—decreasing slack times means changing to faster shipping mode with a shorter delivery window (in order to meet the deadline for orders). Although the actual cost incurred by retailers is confidential, we understand that the actual rates have the same format, but linear discounts are used and these are company dependent.

¹⁸Note that this is an imposed assumption to ensure shipment of orders. In the rest of the paper, unless noted, we only specify the function in the region of $(1, \infty)$.

following shipping decisions: (1) For each order type, which subset of orders should be shipped? (2) From which warehouse should these orders be shipped? (3) How should we split them into packages/delivery tasks? After the delivery, if any, is executed, m new orders of type X arrive with probability $\alpha_{X,m}$ —new orders arrive at the end of each period. For ease of notation, we assume that the arrivals of different order types are independent, e.g., the probability that we have two new orders of type A and three new orders of type B in a period is $\alpha_{A,2} \alpha_{B,3}$.

We now provide the Dynamic Programming (DP) formulation of our problem. For each $s_{X,i}$ in \vec{s} , let $x_{X,i}$ denote the decision whether to ship the corresponding order (i.e., $x_{X,i} = 1$ if the order is shipped and $x_{X,i} = \infty$ otherwise). For order type B , let h_i^W denote the decision whether to ship it from warehouse $W \in \{1, 2\}$ (i.e., if $h_i^W = 1$, it is to be shipped from W ; if $h_i^1 = h_i^2 = \infty$, the order is not shipped; since we only ship from one of the warehouses, the case $h_i^1 = h_i^2 = 1$ is impossible). Let $C_W(\vec{s}')$ denote the minimum cost of shipping a subset \vec{s}' from warehouse W in the current period and let $f(\cdot)$ denote the total shipping costs from both warehouses in the current period (i.e., $f(\cdot) = C_1(\cdot) + C_2(\cdot)$). The rigorous definition of $f(\cdot)$ can be found below in equation (1.1). $C_W(\vec{s}')$ depends on the items shipped and has an “iterative” form since it is not necessary to put all the items of \vec{s}' in a single package.¹⁹ However, it is easy to establish that (see Lemma I.14 in Section 1.6), for a given warehouse, if it is optimal to ship two orders with slack times θ_i and θ_j ($\theta_i \leq \theta_j$) in the same package from the same warehouse, then it is also optimal to ship all orders with slack times between θ_i and θ_j in this package. For any vectors v and $\mu \in \mathbb{R}^n$, we define $v \cdot \mu = (v_1\mu_1, v_2\mu_2, \dots, v_n\mu_n)$ and $\rho(v, \mu)$ to be a function whose output is the ordered list of all non- ∞ elements of both v and μ combined (e.g., if $v = (1, 5, \infty)$ and $\mu = (3, 4, 6)$, then $\rho(v, \mu) = (1, 3, 4, 5, 6)$).

¹⁹Consider a simple case with two orders to ship, one with slack time 1 and the other with slack time 10. Suppose that the costs are given by $F(1) = 20$, $v(1) = 10$, $F(10) = 2$, and $v(10) = 1$. The total costs of shipping the orders in two separate packages is $20+10+2+1 = 33$, which is smaller than the total costs of shipping both orders in one package $20+2 \times 10 = 40$.

Let $V_t(\vec{s})$ denote the cost-to-go function at the beginning of period t . We write $V_t(\cdot)$ recursively as follows:

$$\text{For } t > 1: V_t(\vec{s}) = \min_{\vec{x}} \left\{ f(\vec{s}_A \cdot \vec{x}_A, \vec{s}_B \cdot \vec{x}_B, \vec{s}_C \cdot \vec{x}_C) + \mathbf{E}[V_{t-1}(\vec{s}_{\vec{x}})] \right\} \quad (1.1)$$

$$\text{For } t = 1: V_1(\vec{s}) = f(\vec{s}_A \cdot \vec{x}_A, \vec{s}_B \cdot \vec{x}_B, \vec{s}_C \cdot \vec{x}_C)$$

$$\text{where } f(\vec{y}_A, \vec{y}_B, \vec{y}_C) = \min_{\vec{h}} \{ C_1(\rho(\vec{y}_A, \vec{y}_B \cdot \vec{h}^1)) + C_2(\rho(\vec{y}_B \cdot \vec{h}^2, \vec{y}_C)) \},$$

$$x_{X,i} = \min\{h_i^1, h_i^2\},$$

$$C_W(\theta_1, \theta_2, \dots, \theta_n) = \min \left\{ \begin{array}{l} F_W(\theta_1) + v_W(\theta_1) + C_W(\theta_2, \dots, \theta_n) \\ \quad \text{ship one order in the 1}^{st} \text{ package} \\ F_W(\theta_1) + 2v_W(\theta_1) + C_W(\theta_3, \dots, \theta_n) \\ \quad \text{ship two orders in the 1}^{st} \text{ package} \\ \dots \\ F_W(\theta_1) + nv_W(\theta_1) + C_W(\emptyset) \\ \quad \text{ship } n \text{ orders in the 1}^{st} \text{ package} \end{array} \right.$$

$$C_w(\emptyset) = 0, \quad \forall i, n, W \in \{1, 2\}, \theta_1 \leq \dots \leq \theta_n$$

To ensure that all pending orders with $s_{X,i} = 1$ are shipped, we impose a boundary condition $V_t(\vec{s}) = \infty$ for all \vec{s} with $\min_{i,X} s_{X,i} = 0$. Note that $\vec{s}_{\vec{x}}$ is the new vector of slack times resulting from shipping decisions \vec{x} . We can calculate $\vec{s}_{\vec{x}}$ in two steps. We first construct a vector $\vec{\hat{s}}$ as follows: For all $s_{X,i} < \infty$, if $x_{X,i} = \infty$, let $\hat{s}_{X,i} = s_{X,i} - 1$; if $x_{X,i} = 1$, which means that the corresponding order is shipped, then it is no longer included in $\vec{\hat{s}}$. Next, we can construct $\vec{s}_{\vec{x},X}$ by appending $\vec{\hat{s}}_X$ with a vector of new arrivals of type X , i.e., $\vec{s}_{\vec{x},X} = \{\vec{\hat{s}}_X, \underbrace{d, \dots, d}_m\}$ with probability $\alpha_{X,m}$. Observe that the cost function $C_i(\cdot)$ incorporates the optimal splitting of the shipment into several packages. It does not have an explicit mathematical expression and can only

be captured recursively, making the above DP difficult to analyze.²⁰ In the following sections, we analyze three simplified cases to get a better understanding of the structure of the optimal policy.

1.4 Single Warehouse with Only Fixed Cost

We first study the simplest case where all orders can be shipped from a single warehouse, i.e., all orders are of the same type. Suppose that there are currently n pending orders. While there are $2^n - 1$ different ways of choosing which orders to ship, it is not difficult to show that it is always optimal to either ship all n orders at the same time or none at all. Moreover, all orders are shipped in a single package, i.e., no package splitting takes place.

Lemma I.1. *(a) If it is optimal to ship at least one order in the current period, then it is optimal to ship all pending orders in the same period. (b) All orders are shipped in one package. (c) The total shipping costs for this package is a function of the smallest slack time among the shipped orders.*

Lemma I.1 implies state-dimensionality reduction, i.e., it allows us to simply use the smallest slack time of all pending orders (i.e., $z := \min\{s_1, s_2, \dots\}$), instead of all the slack times (\vec{s}), as the state variable. Let $\alpha = \sum_{m=1}^{\bar{m}} \alpha_m$, where \bar{m} is the upper

²⁰Note that the function C_w is high-dimensional and is neither convex nor concave. Thus, it is difficult to obtain any structural properties. However, in practical settings, the number of comparison is reasonably small and it is easy to solve C_w numerically. Thus, the retailers can easily calculate the package splitting decision once a shipping policy is given.

bound of number of new orders. The DP formulation in (1.1) can be simplified into:

For $t > 1$ and $1 \leq z \leq d$, we have:

$$V_t(z) = \min \begin{cases} F(z) + V_t(\infty) & \text{Ship} \\ V_{t-1}(z-1) & \text{Do not ship} \end{cases} \quad (1.2)$$

$$V_t(\infty) = \alpha V_{t-1}(d) + (1 - \alpha)V_{t-1}(\infty). \quad (1.3)$$

For $t = 1$, we have: $V_1(z) = F(z)$.

Shipping all pending orders in period t incurs a current cost $F(z)$ plus a future cost $V_t(\infty)$, while holding all orders to the next period reduces the slack time by one. Equation (1.3) corresponds to the case where there is no pending order, i.e., either new orders arrive with probability α and the slack time becomes d , or no order arrives and the slack time remains ∞ . The above formulation implies that all orders must be shipped by the end of the horizon. The following proposition describes a property of $V_t(\cdot)$.

Proposition I.2. $V_t(z)$ is non-increasing in $z \geq 1$ given t and is non-decreasing in t given $z \geq 1$.

Proposition I.2 has an intuitive interpretation: Smaller slack time means more urgency, which implies higher expected total shipping costs; smaller t means fewer future orders, which implies smaller expected total shipping costs.

1.4.1 The Optimal Policy: Its Structure and Properties

We now show that the optimal shipping policy has a simple threshold structure. We first state a lemma that will be used to prove this property.

Lemma I.3. *For all $t \geq 1$ and $z \geq 2$, $V_t(z-1) - V_t(z) \geq F(z-1) - F(z)$.*

Lemma I.3 means that if the critical order, the order with smallest slack time, becomes more urgent, then the impact on the cost-to-go function is *larger* than that on the current shipping costs (i.e., if it was shipped in the current period). This result may seem counter-intuitive: One might argue that the difference between two optimal values, one for slack of $z - 1$ and another for slack of z should be smaller than, or equal to, the difference between the corresponding costs in the current period. This would be because the optimal value functions have more flexibility and allow for many alternatives which may make the difference between options that start with $z - 1$ and z smaller (closer to each other). However, this intuition is incorrect in this case. $V_t(z-1)$ may correspond to the optimal next shipment taking place in some period t_1 whereas $V_t(z)$ may correspond to the optimal next shipment taking place in a different period t_2 , where both t_1 and t_2 could be different from the current period t . Therefore, without additional assumptions, either the left hand side or right hand side can be bigger.²¹ Lemma I.3 provides a link between the cost-to-go functions and the current shipping costs. It allows us to explicitly compare the shipping costs of different alternatives.

Theorem I.4. *There exists an integer threshold τ_t such that the optimal decision in period t is to hold all pending orders if $z > \tau_t$ and to ship all of them if $z \leq \tau_t$.*

PROOF. Fix time period t . To prove the existence of a threshold τ_t , it is sufficient to show that if the optimal decision for slack time $z \geq 3$ is to ship all orders, then the optimal decision for slack time $z - 1$ is also to ship all orders. (The case $z = 2$ and $z = 1$ are trivial because we must ship when slack time becomes 1.) By DP

²¹One can express the difference between two optimal values as a comparison between the minimum cost in two sets of situations, which start with $z - 1$ or with z , and include the alternatives that ship with the same delay of k periods from the current period. The difference between shipping costs k period later, $F(z - k - 1)$ and $F(z - k)$, increases with k (due to the convexity of shipping costs), which drives the increases in the differences of the cost-to-go function in this problem. The result can be viewed as a natural extension of a simple mathematics property that the difference between the minimum of two sets of values is larger than the minimum difference between each element in these two sets.

formulation, it is optimal to ship if, for slack time z , $F(z) + V_t(\infty) \leq V_{t-1}(z - 1)$. Then, for slack time $z - 1$, $F(z - 1) + V_t(\infty) = (F(z - 1) - F(z)) + (F(z) + V_t(\infty)) \leq (F(z - 2) - F(z - 1)) + V_{t-1}(z - 1) \leq V_{t-1}(z - 2)$, where the last inequality follows by Lemma I.3. This implies that it is also optimal to ship for slack time $z - 1$, which completes the proof.

Theorem I.4 shows the existence of threshold τ_t for each time t . Since we assume a stationary arrival probability, using the standard convergence argument as in the infinite-horizon literature (Gosavi, 2003), it is not difficult to show that there exists some τ^* , such that $\tau_t \rightarrow \tau^*$ as $t \rightarrow \infty$. The following theorem tells us that the optimal constant threshold is easy to compute.

Theorem I.5. *Suppose we use a constant threshold τ in all periods. Then, the expected average shipping cost during T periods converges to $G_\tau(\alpha, d) = F(\tau) \left(\frac{1}{\alpha} + d - \tau\right)^{-1}$ as $T \rightarrow \infty$.*

The interpretation is straightforward. $F(\tau)$ is the shipping cost incurred when the pending orders trigger the threshold τ . $\left(\frac{1}{\alpha} + d - \tau\right)^{-1}$ is the average cycle time: after all pending orders are shipped in the last cycle, it takes $\frac{1}{\alpha}$ on average until a new order arrives and $d - \tau$ periods until all new orders are shipped. To calculate the optimal τ^* , we simply need to minimize $G_\tau(\alpha, d)$ over the set $\{1, 2, \dots, d\}$. One simple application of Theorem I.5 is for the case where $F(\cdot)$ is linear:

Proposition I.6. If the fixed cost is linear in slack times, then the optimal threshold τ is either 1 or d .

The following lemma describes the behavior of τ^* as a function of α .

Lemma I.7. *The optimal constant threshold τ^* is non-increasing as α increases.*

Lemma I.7 implies that τ^* is smallest when $\alpha \approx 1$ and is largest when $\alpha \approx 0$. This is quite intuitive: If orders arrive very frequently, the opportunity to consolidate orders

is high, which provides an incentive for the retailer to delay shipping. If, on the other hand, orders arrive infrequently, the opportunity to consolidate orders is low. Thus, it is better to ship the pending orders earlier due to the risk of incurring higher shipping costs without the benefit of consolidation.

1.5 Two Warehouses with Only Fixed Cost

We now consider the case where the retailer fulfills orders from two warehouses and orders can be classified into three types as in Section 1.3. The cost structure, similar to Section 1.4, includes only the fixed cost. Lemma I.8 below is similar to Lemma I.1.

Lemma I.8. (a) *If it is optimal to ship an order of a particular type in the current period, then it is optimal to ship all pending orders of the same type.* (b) *The incurred shipping cost in a period is a function of the vector of the smallest slack times of respective types among the shipped orders.*

By Lemma I.8, the vector of the smallest slack times (z_A, z_B, z_C) , for orders of types A, B , and C , respectively, completely describes the state space. Let $\alpha_X = \sum_{m \in N} \alpha_{X,m}$. The DP formulation in (1.1) can be simplified to:

$$\text{For } t > 1: V_t(z_A, z_B, z_C) = \min_{(x_A, x_B, x_C)} \{f(z_A x_A, z_B x_B, z_C x_C) + \mathbf{E}[V_{t-1}(\tilde{z}_A, \tilde{z}_B, \tilde{z}_C)]\}$$

$$\text{For } t = 1: V_1(z_A, z_B, z_C) = f(z_A, z_B, z_C)$$

$$\text{where } f(y_1, y_2, y_3) = \min\{F_1(\min\{y_1, y_2\}) + F_2(y_3), F_1(y_1) + F_2(\min\{y_2, y_3\})\}, \forall y_1, y_2, y_3$$

with boundary conditions $V_t(0, \cdot, \cdot) = V_t(\cdot, 0, \cdot) = V_t(\cdot, \cdot, 0) = \infty$, and \tilde{z} denoting the new vector of slack times resulting from shipping decisions (x_A, x_B, x_C) . Similar to Proposition I.2, we have:

Proposition I.9. Suppose that $1 \leq z'_A \leq z_A, 1 \leq z'_B \leq z_B$, and $1 \leq z'_C \leq z_C$. For all $t \geq 1$, we have: $V_t(z'_A, z'_B, z'_C) \geq V_t(z_A, z_B, z_C)$ and $V_t(z_A, z_B, z_C) \geq V_{t-1}(z_A, z_B, z_C)$.

1.5.1 The Optimal Policy: Its Structure and Properties

We first describe the optimal policy for the case where W1 and W2 are symmetric, i.e., $F_1(\cdot) = F_2(\cdot)$ and $\alpha_A = \alpha_C$. We then briefly discuss the complexity of the optimal policy in the general case. The following two lemmas are useful for the analysis.

Lemma I.10. *Suppose that $z_A, z_B < \infty$. If it is optimal to ship orders of type B from W1 in the current period, then it is also optimal to ship orders of type A from W1 in the same period. By symmetry, if $z_B, z_C < \infty$ and it is optimal to ship orders of type B from W2 in the current period, then it is also optimal to ship orders of type C from W2 in the same period.*

Lemma I.11. *For all $t \geq 1$ and $z_A, z_B, z_C \geq 2$, the following holds:*

1. *If $z_A \leq z_B$, then $V_t(z_A - 1, z_B, z_C) - V_t(z_A, z_B, z_C) \geq F_1(z_A - 1) - F_1(z_A)$.*
2. *If $z_C \leq z_B$, then $V_t(z_A, z_B, z_C - 1) - V_t(z_A, z_B, z_C) \geq F_2(z_C - 1) - F_2(z_C)$.*
3. *If $z_B \leq \min\{z_A, z_C\}$, then $V_t(z_A, z_B - 1, z_C) - V_t(z_A, z_B, z_C) \geq \min\{F_1(z_B - 1) - F_1(z_B), F_2(z_B - 1) - F_2(z_B)\}$.*

Lemma I.10 simplifies the shipping alternatives by eliminating the possibility of shipping B alone from either W1 or W2,²² and Lemma I.11 is the analog of Lemma I.3. The formal definition of the optimal policy and the corresponding six boundaries are given below.

²²It is worth noting that “shipping only A from W1” can still be optimal. A simple example is when $z_C < z_B < z_A$ and order A is shipped in the current period, but it is more efficient to ship B later with C. A more detailed example (using the cost structure of Breakaway Courier) can be found in the Introduction Section.

Theorem I.12. *In the symmetric two-warehouse setting, for $z_A, z_B, z_C \geq 1$, there exist six boundaries $\tau_{A,t}^{AB}(z_B, z_C) \leq \tau_{A,t}^A(z_B, z_C)$, $\tau_{C,t}^{BC}(z_A, z_B) \leq \tau_{C,t}^C(z_A, z_B)$, $\tau_{B,t}^1(z_A, z_C)$ and $\tau_{B,t}^2(z_A, z_C)$ that completely characterize an optimal policy in the symmetric two-warehouse setting. For any given state (z_A, z_B, z_C) , if z_X , $X \in \{A, B, C\}$, crosses the boundary, then it is optimal to ship orders of type X , following the policy below.*

1. *If $z_A \leq \tau_{A,t}^{AB}(z_B, z_C)$, it is optimal to ship both orders of type A and B from W1;
if $z_A \leq \tau_{A,t}^A(z_B, z_C)$, it is optimal to ship orders of type A from W1;*
2. *If $z_B \leq \tau_{B,t}^1(z_A, z_C)$, it is optimal to ship both orders of type A and B from W1;
if $z_B \leq \tau_{B,t}^2(z_A, z_C)$, it is optimal to ship orders of type B and C from W2;*
3. *If $z_C \leq \tau_{C,t}^{BC}(z_A, z_B)$, it is optimal to ship both orders of types B and C from W2;
if $z_C \leq \tau_{C,t}^C(z_A, z_B)$, it is optimal to ship orders of type C from W2.*

Moreover, the following also hold: $\tau_{A,t}^{AB}(\infty, z_C) = \tau_{A,t}^A(\infty, z_C)$ and $\tau_{C,t}^{BC}(z_A, \infty) = \tau_{C,t}^C(z_A, \infty)$.

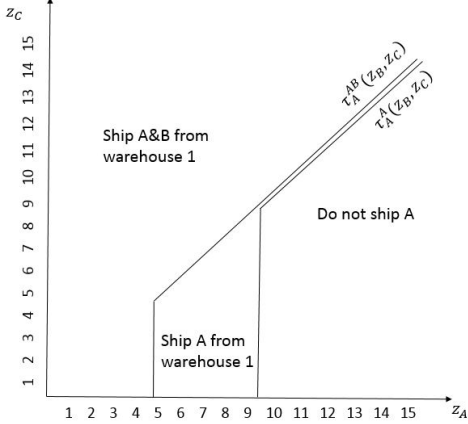


Figure 1.1: Boundaries for type A orders when $z_B = 5$

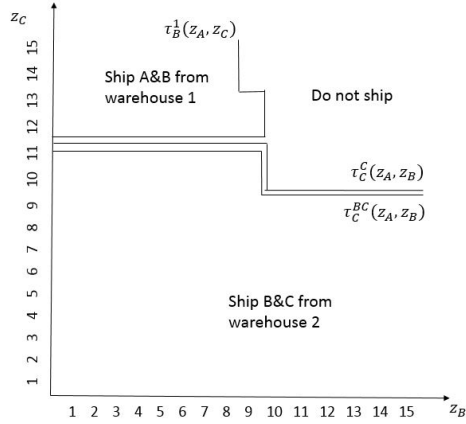


Figure 1.2: Boundary τ_B^1 for type B orders when $z_A = 11$

Each of the boundaries in the above lemma is a function of two slack variables (with either slack variable z_A or z_B fixed), so they can be viewed as surfaces in three-dimensional space. These boundaries completely characterize the optimal shipping

policy.²³ To visualize the six boundaries in Theorem I.12, we provide two figures which show 2-dimensional “cuts” from the optimal-decision 3-dimensional space. Figure 1.1 provides an illustration of the boundaries of type A order, and Figure 1.2 provides an illustration of one of the boundaries of type B orders, $\tau_B^1(z_A, z_C)$. We omit the illustrations for boundaries of type C orders and $\tau_B^2(z_A, z_C)$ of type B orders. Obviously, the 2-dimensional cuts provide incomplete information. For example, in Figure 1.1, the position of z_A and z_C describes an action with respect to order of type A, but still does not fully determine the action for orders of type B (which only becomes clear when one looks at z_B).²⁴

In some cases, the six boundaries in the two-warehouse setting can be reduced to thresholds, similar to the structure in one-warehouse setting. In an extreme case, when $\alpha_B = 0$, since orders of type B do not exist, W1 is independent of W2 and the optimal policy for orders types A and C are each characterized by time-dependent thresholds. This can also be observed based on Theorem I.12: Since $\alpha_B = 0$, the slack time of orders of type B always equals ∞ . As $\tau_{A,t}^{AB}(\infty, z_C) = \tau_{A,t}^A(\infty, z_C)$ and $\tau_{C,t}^{BC}(z_A, \infty) = \tau_{C,t}^C(z_A, \infty)$, the six boundaries reduce to only two boundaries. In general, however, all six boundaries are required to properly define the optimal shipping policy. The following lemma is the analog of Lemma I.7.

Lemma I.13. *The optimal stationary boundaries are all decreasing in α_X , $X \in \{A, B, C\}$.*

Note that when the boundaries of one type of orders change, there is a chain effect

²³Note that the optimal policy may be just “shipping only A from W1” or “shipping both B and C from W2”. When $z_A > z_B$, shipping B with A increases the shipping cost in the current period, and shipping B with C (in the current or later period) may be cheaper.

²⁴Parameters: $F_1(z) = F_2(z) = -0.005z^3 + 0.7z^2 - 16z + 109.6$, $d = 15$, $\alpha = (0.55, 0.2, 0.55)$. Note that in both figures, we have states that are indifferent between shipping alternatives, e.g., in Figure 1.1, there are ties of shipping alternatives for the region of $z_A \leq 4$ and $z_C \leq 4$ – the costs are equal between “shipping B with A” and “shipping B with C.” Similarly, in Figure 1.2, the line where three boundaries coincide for $z_C = 11$ is also the line of indifference.

on orders of other types: The intuition is that increasing the arrival probability for order type A leads to a lower boundary for order type B , which provides more opportunities for orders type C to be jointly shipped with orders type B without incurring additional cost.

We now discuss the case where $W1$ and $W2$ are asymmetric, i.e., either the cost functions of the two warehouses or the arrival probability of order types A and C are not the same. In such a case, the optimal policy can no longer be characterized by the six boundaries in Theorem I.12, see Figures 1.3 and 1.4 for illustrations (The optimal policy is calculated using the dynamic programming (1.1)). In Figure 1.3, for a fixed z_A , when z_B decreases, the optimal solution for state (z_A, z_B, z_C) can change from “Do not ship” to “ship A and B from warehouse 1” and then to “Do not ship” again. In Figure 1.4, for a fixed z_C , the region of “Ship A from warehouse 1” is not even connected. Since the optimal policy for the asymmetric case can be very complex, we do not study its structural properties; instead, we propose simple heuristics that can perform well for most cases. We discuss them next.



Figure 1.3: Asymmetric case 1 when $z_A = 8$

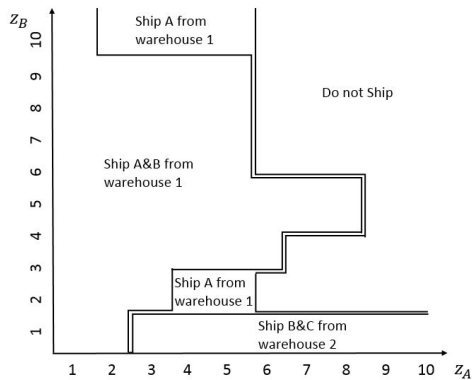


Figure 1.4: Asymmetric case 2 when $z_C = 9$

1.5.2 Simple Heuristics

We now propose two heuristics: *warehouse-based* heuristic, which utilizes two thresholds, one for each warehouse, and *order-based* heuristic, which utilizes three thresholds,

one for each order type. Their performances are evaluated using percentage gap, defined as $(C_H - C^*)/C^*$, where C_H is the expected average costs per period by applying the heuristic and C^* is the expected average costs under the optimal policy. The two heuristics perform very similarly.

1.5.2.1 Warehouse-based Heuristic.

The warehouse-based heuristic assigns a constant threshold to each warehouse (τ_1 for $W1$ and τ_2 for $W2$). If the slack time of an order type falls below the threshold for the corresponding warehouse, then all pending orders for that warehouse are shipped. Note that this heuristic is not equivalent to treating the two warehouses independently as it allows orders type B to be fulfilled from either warehouses. The formal description of the heuristic follows:

Warehouse-Based Heuristic

Given (τ_1, τ_2) and (z_A, z_B, z_C) , do:

1. If $\min\{z_A, z_B\} \leq \tau_1$, ship orders of types A and B from $W1$ and update $z_A = z_B = \infty$;
2. If $\min\{z_B, z_C\} \leq \tau_2$, ship orders of types B and C from $W2$ and update $z_B = z_C = \infty$.

The warehouse-based heuristic is easy to implement²⁵. It simplifies the decisions by bundling the shipment of orders types A and B (and C and B) together. The thresholds (τ_1, τ_2) need to be optimized using simulation-based optimization. In the numerical experiments, we do this by running a complete search over all possible values

²⁵This type of warehouse-based operations is widely used in practice (Xu et al., 2009; Acimovic and Graves, 2015). Such warehouse-based structure is also consistent with what we learned from the retailer in China.

of the thresholds and, for each (τ_1, τ_2) , we estimate the corresponding C_H using 100 Monte Carlo simulations.

1.5.2.2 Order-based Heuristic.

The order-based heuristic assigns a constant threshold to each order type (τ_X for order type $X \in \{A, B, C\}$). If the threshold of a particular order type is triggered, all pending orders of this type are shipped and orders of other types may also be shipped jointly according to a pre-specified consolidation rule. For each order type, this heuristic essentially replaces the two boundaries in the optimal policy with one constant threshold. Further, it uses a myopic consolidation rule to decide whether to ship the order type alone, or together with other orders from the same warehouse. The order-based heuristic is more flexible than the warehouse-based heuristic as it allows shipping order types A and B (or B and C) together or separately depending on the realization of orders. That said, the order-based heuristic is also more complex. In the order-based heuristic, when the threshold of an order type is triggered, the retailer needs to properly decide whether to consolidate this order with other order types.

In the order-based heuristic, we use a *one-period myopic* consolidation rule. Under our proposed rule, consolidation is decided by finding the best alternative, as if all orders had to be shipped in the current period. The details are shown below.

Order-Based Heuristic

Given (τ_A, τ_B, τ_C) and (z_A, z_B, z_C) , do:

1. If $z_A \leq \tau_A$: if $F_1(\min\{z_A, z_B\}) + F_2(z_C) \leq F_1(z_A) + F_2(\min\{z_B, z_C\})$,²⁶ ship all orders of types A and B from W1; otherwise, ship only all orders of type A (no consolidation);

²⁶These are formulations of the one-period myopic cost to ship B from either W1 or W2

2. If $z_B \leq \tau_B$: if $F_1(\min\{z_A, z_B\}) + F_2(z_C) \leq F_1(z_A) + F_2(\min\{z_B, z_C\})$, ship orders of types A and B from W1; otherwise, ship all orders of types B and C from W2;
3. If $z_C \leq \tau_C$: if $F_1(z_A) + F_2(\min\{z_B, z_C\}) \leq F_1(\min\{z_A, z_B\}) + F_2(z_C)$, ship all orders of types B and C from W2; otherwise, ship only all orders of type C (no consolidation);
4. If more than one threshold is crossed at the same time, orders of type B are the highest priority and then A and C (e.g., if τ_B and τ_C are crossed, we first proceed according to point 2 and then to point 3).

Similar to the warehouse-based heuristic, the thresholds (τ_A, τ_B, τ_C) need to be optimized. In the experiments, we do this by running a complete search and, for each (τ_A, τ_B, τ_C) , we estimate C_H using 100 Monte Carlo simulations.

1.5.2.3 Simulation results.

The performance of both the warehouse-based and the order-based heuristics are tested using simulations. Unless otherwise noted, we use d equals 5. This can be interpreted either as retailers have five shipment options to deliver orders, or the deadline of orders is five periods. Both interpretations are close to the current logistics practice. The customer-ordering frequency for each order type $(\alpha_x, x \in \{A, B, C\})$ varies within $(0,1)$, capturing a wide range of customers from those who purchase very infrequently to those who purchase very frequently. The cost function is based on UPS rates, in terms of both the functional form and the parameters.²⁷ Both symmetric and asymmetric settings are tested. The asymmetric setting captures the potential asymmetry

²⁷Although the actual cost incurred by retailers is confidential due to discounts, we believe that the structure is similar to the official (published) rate of logistics firms and that their structures are not dramatically different from one to another. We average several UPS rates, and since they have a convex shape, we fit it with a quadratic function.

in the order frequency, $\alpha_A \neq \alpha_C$, and the asymmetry in the cost structures of the two warehouses: $F_1(\cdot) = F_2(\cdot) + \beta_1$ or $F_2(\cdot) = \beta_2 F_1(\cdot)$. These two forms of cost functions capture the asymmetry driven by either the difference of distance from the warehouse to the customer (the former additive function), or the various cost “efficiency” of shipping orders from different warehouses (the latter multiplicative function). β_1 and β_2 vary from low to high across a wide range. This results in 100 cases capturing different economics and demand frequencies. Details of the simulation can be found in appendix ???. Unless otherwise noted, the tests are run on a PC with 3.5GHz CPU and 16GB memory.

For the warehouse-based heuristic, the total running time to compute the best (τ_1, τ_2) is only 1.79 seconds CPU time on average.²⁸ The warehouse-based heuristic performs very well with the average percentage gap of 1.93% and a standard deviation of 2.39%. This suggests that the six boundaries of the optimal policy can be replaced with two simple warehouse thresholds. For the order-based heuristic, the total running time to compute the best (τ_A, τ_B, τ_C) is 6.90 seconds CPU time on average. Our numerical experiments show that the order-based heuristic performs very well, with the average percentage gap of 0.11% and a standard deviation of 0.25%. Surprisingly, this is despite the fact that the one-period myopic consolidation rule ignores expected total costs for future periods. We conclude that both warehouse-based and order-based heuristic policies perform sufficiently well. Their performance suggest that the six optimal boundaries in Theorem I.12 can be well-approximated by constant thresholds, either corresponding to warehouses or order types. These approximations not only provide easy-to-implement and intuitive heuristics, but also justify the use of these heuristics in a more complex setting in Section 1.7.

²⁸Note that the values of thresholds are pre-calculated before the heuristics are implemented by retailers. Given the thresholds, a retailer directly executes the shipping policy, without any need for further optimization. Thus, the CPU time is an up-front investment and it does not influence the implementation of the heuristics.

1.6 One Warehouse with Fixed and Variable Costs

In Section 1.5, we extend the analysis of the one-warehouse case with only fixed cost to the two-warehouse setting. We now extend the analysis of one-warehouse case to the setting with both fixed and variable costs. While the optimal policy for the case with only fixed cost is easy to characterize, the optimal policy for the case with both fixed and variable costs is more challenging, even for the setting with only one warehouse. First, it may no longer be optimal to ship all pending orders in one package (i.e., package splitting is possible), as discussed in Section 1.3. Second, when some orders are shipped, it may no longer be optimal to ship all pending orders in the same period (in contrast to Lemma I.1). Indeed, it may now be more economical to intentionally delay the shipment of some orders to consolidate them with future orders.²⁹ In what follows, we first discuss the shipping costs of orders and the simplified DP formulation. Next we show, for some cases, that the optimal policy is a volume-dependent threshold. To further simplify this policy, we approximate the optimal threshold with a constant and show, using numerical experiments, that this approximation incurs a very small additional cost compared to the optimal policy.

1.6.1 Dynamic Programming (DP) Formulation

To evaluate the shipping costs of orders, we first need to consider how to split them into shipments. Given a list S of n orders to be shipped, with slack times $z_{s_1}, z_{s_2}, \dots, z_{s_n}$ ($z_{s_1} \leq z_{s_2} \leq \dots \leq z_{s_n}$), the number of all possible package-splitting alternatives is very high (i.e., $\frac{1}{e} \sum_{k=0}^{\infty} \frac{k^n}{k!}$, according to Lovasz, 1993). Fortunately, under the optimal policy,

²⁹Consider the case where $d = 10$ and there are two orders to ship, with slack time 1 and 10, respectively. The costs are $F(1) = 20$, $v(1) = 10$, $F(9) = 2.5$, $v(9) = 1.25$, and $F(10) = 2$, $v(10) = 1$. Anticipating that a new (third) order arrives in the next period, the cost of shipping only the order with slack time of 1 in the current period and shipping orders 2 and 3 in the next period is $20 + 10 + 2.5 + 2 \times 1.25 = 35$, which is smaller than the cost of shipping orders 1 and 2 in the current period and shipping the new order 3 in the next period $20 + 10 + 2 + 1 + 2 + 1 = 36$.

the optimal splitting and shipping is *continuous* in the sense described below in Lemma I.14.

Lemma I.14. *For any state, if, in the optimal solution, z_i is the most urgent order to be shipped and n orders are to be shipped **in the current period**, then the other orders shipped in the current period should be those corresponding to slack times $z_{i+1}, \dots, z_{i+n-1}$.*

*For orders in a set S to be shipped, if, in the optimal solution, z_{s_i} is the most urgent order in a package and n orders are to be shipped **in this package**, then the other orders included in the same package should be those corresponding to slack times $z_{s_{i+1}}, \dots, z_{s_{i+n-1}}$.*

The optimal shipping and splitting policy above is referred to as continuous. Note, however, that even when the shipping policy is continuous, it does not imply that most urgent orders are shipped. For example, for state $(z_1, z_2, z_3, z_4, z_5, z_6) = (4, 8, 9, 10, 11, 12)$, shipping orders $(8, 9, 10, 11, 12)$ and leaving order (4) to future periods can be more economical compared to shipping orders $(4, 8, 9, 10, 11)$.³⁰ The lack of clarity on which order to ship together with the lack of closed-form expression for the cost function $C(\cdot)$, see Section 1.3, are the key reasons why the optimal policy for the one-warehouse setting with both fixed and variables costs is difficult to analyze in general. That said, we are able to derive some results for the case when $d \leq 3$ and at most one order arrives in each period, i.e., $\alpha_m = 0$ for $m \geq 2$ and $\alpha = \alpha_1$. We state these in the following theorem.

Theorem I.15. *Suppose that $d \leq 3$ and at most one order arrives in each period.*

Then, the optimal policy has the following properties:

³⁰Consider the case where $d = 12$, $F(\cdot) = 15$, $v(z) = 13 - z$ for $0 < z \leq 8$ and $v(z) = 5$ for $8 < z \leq 12$. In period t , for state $(z_1, z_2, z_3, z_4, z_5, z_6) = (4, 8, 9, 10, 11, 12)$, suppose a policy suggests to ship five orders in period t and the remaining order with a new-arrival order z_7 in period $t - 3$. Shipping orders $(z_1, z_2, z_3, z_4, z_5)$ in period t and (z_6, z_7) in period $t - 3$ incurs cost of 85, while shipping $(z_2, z_3, z_4, z_5, z_6)$ in period t and (z_1, z_7) in period $t - 3$ incurs less cost, 79.

1. *It is optimal to either ship all pending orders or not ship any order at all at any state that is reachable³¹ under the optimal policy.*
2. *Given n pending orders, the optimal policy is characterized by volume-dependent threshold $\tau_t(n)$. Specifically, if $z_1 \leq \tau_t(n)$, it is optimal to ship all pending orders; otherwise, it is optimal not to ship any order.*
3. *For all t , the threshold $\tau_t(n)$ is non-decreasing in n .*

Per Theorem I.15, instead of having to consider how many orders to ship, the retailer only needs to choose between shipping all orders and shipping no order. The only caveat is that the threshold is a function of the volume of pending orders. The intuition is that, with a large number of pending orders in the system, the cost of waiting for future orders increases, as there are more orders that are potentially affected by the change to more expensive shipping modes. This decreases the incentive to wait for future orders, which is equivalent to a higher threshold. It is worth noting that, when it is optimal to ship all pending orders, it may still be optimal to ship orders in more than one package.³² That is, although the shipping policy is decided only by the smallest slack time z_1 and the volume of pending orders m , it is still necessary to keep track the slack times of all pending orders as they affect both the package-splitting and the total shipping costs. While the results of Theorem I.15 are only proved for the case where $d \leq 3$ and at most one arrival per period (due to analytical tractability), our numerical tests across different instances where d varies from 6 to 10 and multiple

³¹Reachable states are defined for a given policy. Assume that the initial state is (∞) , with no orders in the system. Then, a state is reachable if there exists a sequence of demand arrivals such that the system will move into this state, while following the given policy. Not all states are reachable. Consider a very simple example where we ship each order when the slack time is 2. For that policy no state with slack = 1 is reachable.

³²Consider the case where the orders have two possible volumes, 1 or 100, with arrival probability α_1 and α_2 respectively. The deadline of orders is $d = 2$. The cost function is $F(1) = 100$, $F(2) = 99$, $v(1) = 2$, $v(2) = 1$. In any period, for state where there are two orders, one with volume 1 and slack time 1, the other with volume 100 and slack time 2, it is optimal to ship both orders but in two separate packages.

orders can arrive in one period show that these results continue to hold in absolutely all cases. Assuming the structures of the policy in Theorem I.15, the DP formulation (1.1) can be simplified into the following:

$$\text{For } t > 1, V_t(z_1, z_2, \dots, z_n) = \min \begin{cases} C(z_1, z_2, \dots, z_n) + V_t(\infty) & \text{Ship } m \text{ orders} \\ \alpha V_{t-1}(z_1 - 1, z_2 - 1, \dots, z_n - 1, d) \\ \quad + (1 - \alpha)V_{t-1}(z_1 - 1, z_2 - 1, \dots, z_n - 1) \\ \text{Do not ship} \end{cases}$$

$$V_t(\infty) = \alpha V_{t-1}(d) + (1 - \alpha)V_{t-1}(\infty).$$

$$\text{For } t = 1, V_1(z_1, z_2, \dots, z_n) = C(z_1, z_2, \dots, z_n).$$

1.6.2 A Simple Heuristic

As the optimal policy in Theorem I.15 can be characterized by volume-dependent thresholds and the package splitting decision requires another layer of optimization, the policy is still not easy to implement.³³ Motivated by the constant-threshold heuristic developed in Section 1.4, we propose to (1) replace the volume-dependent threshold $\tau_t(n)$ with a constant τ independent of volume of pending orders and (2) only allow orders to be shipped in one package (i.e., no package splitting). The heuristic is formally defined as follows:

Constant-Threshold One-Package Heuristic

Given τ and (z_1, z_2, \dots, z_n) , do:

1. If $z_1 \leq \tau$, ship all n pending orders in one package incurring cost $F(z_1) + nv(z_1)$.
2. Otherwise, wait for future orders.

Similar to the heuristic in Section 1.5, the threshold τ needs to be optimized and

³³Computationally, for $d = 5$, calculating the optimal policy for a one-warehouse setting with fixed and variable costs takes more than 100 times that of a one-warehouse setting with fixed cost.

it can easily be done using complete search and Monte Carlo simulations.

The constant threshold one-package heuristic we propose is labeled as (H3). We conduct numerical experiments to test the performance of our heuristic in the setting where d varies from 4 to 10. The customer-order frequency varies within $(0,1)$. Two cases are tested: (1) at most one arrival per period and (2) up to three orders per period.³⁴ We adjust the probability of order arrivals such that the expected number of orders per period are the same in both cases. The variable costs are set as a fraction of the fixed costs ($v(\cdot) = \gamma F(\cdot)$), where γ varies within $[0.1,1]$, reflecting a wide range of cases from where the fixed cost dominates to where the variable costs dominate. Altogether, this results in 630 cases. Details of simulations can be found in Appendix ??.

For the settings with at most one arrival per period, the percentage gaps compared to the optimal cost are shown in Table 1.1.³⁵ Heuristic H3 is compared with two other heuristics: the constant-threshold heuristic where the thresholds are still approximated by a constant but package splitting is allowed and optimized (H1) and the one-package heuristic where all orders must be shipped in one package and the thresholds can vary with number of pending orders (H2). Note that although allowing volume-dependent thresholds improves the performance of our proposed heuristic, the improvement is very small. Allowing package splitting also leads to extremely small/negligible improvement.³⁶ These results suggest that the optimal policy in the case with fixed and variable costs can be well-approximated by our proposed heuristic H3.

³⁴In our settings, the length of a period naturally corresponds to a delivery option, which could be either several hours or a small number of days. For this length of a period, we have some evidence from the co-founder of one of the largest online retailers in China that the majority of customers place 0 to 3 orders within the actionable time range.

³⁵For the settings with up to three orders per period, the results are very similar—The biggest difference is 0.2% when compared to the results in Table 1.1. We omit the details due to page limit.

³⁶The intuition for package splitting being less critical than volume-based thresholds is as follows: Basically, the correct thresholds decide when all the orders are shipped and they affect the cost of all orders, while the correct package splitting strategy only affects some orders in some packages. Since most of the time packages are not split, imposing non-splitting rule has a very small negative effect.

Heuristic	(H1)	(H2)	(H3)
Varying Thresholds	No	Yes	No
Splitting-package	Yes	No	No
Average	0.24%	0.08%	0.24%
Minimum	0.00%	0.00%	0.00%
10th percentile	0.00%	0.00%	0.00%
25th percentile	0.00%	0.00%	0.00%
50th percentile	0.10%	0.00%	0.10%
75th percentile	0.39%	0.07%	0.39%
90th percentile	0.74%	0.30%	0.74%
Maximum	1.24%	0.97%	1.24%

Table 1.1: Percentage Gaps to Optimal Cost

The good performance of heuristic H3 can also be theoretically justified. Theorem I.17 focuses on large case, and Theorem I.16 focuses on small case of $d \leq 3$ where more convergence properties can be observed. Specifically, Theorem I.17 directly bounds the performance gap between the optimal policy and the constant-threshold one-package heuristic (H3), while Theorem I.16 bounds the performance gap between policies that use varying thresholds (H2) and policies that use constant thresholds (H3). Since H2 performs very close to the optimal policy, Theorem I.16 provides an additional theoretical support for the good performance of the constant-threshold one-package heuristic that we propose.

Theorem I.16. *Suppose that $d \leq 3$, at most one order arrives in each period, and $v(\cdot) = \gamma F(\cdot)$ for some $\gamma \geq 0$. Let \tilde{C} denote the average cost per period for a problem with T periods under the Constant-Threshold One-Package Heuristic (H3) and let C_0 denote the average cost per period under the optimal policy with no package splitting but with varying thresholds. Then, we can bound:*

$$\frac{\tilde{C} - C_0}{C_0} \leq \min \left\{ \frac{\Delta}{a_1 \gamma \Delta + (a_1 + a_2)v(2)}, \frac{F(2)^{\frac{1+(1+3\alpha)(3-\alpha)}{1+2\alpha}}}{F(2)^{\frac{1}{\alpha(1-\alpha)}} + \Delta \frac{1}{\alpha}} \right\}$$

where $\Delta = F(1) - F(2)$, $a_1 = \frac{2}{\alpha} + 1$ and $a_2 = \frac{3}{1-\alpha}$.

Note that the above bound converges to 0 when either $\Delta \rightarrow 0$, $\Delta \rightarrow \infty$, $\alpha \rightarrow 1$, or $\alpha \rightarrow 0$. For the parameters and cost structure used in practice (and also in our simulations), the value of the above bound is at most around 0.09.³⁷

As for larger problems of $d > 3$, we provide a performance bound that directly compares our heuristic (H3) with the optimal policy.

Theorem I.17. *Suppose that $v(\cdot) = \gamma F(\cdot)$ for some $\gamma \geq 0$. Let \tilde{C} denote the average cost per period for a problem with T periods under the Constant-Threshold One-Package Heuristic (H3) and let C^* denote the optimal average cost per period. Then,*

$$\frac{\tilde{C} - C^*}{C^*} \leq \min \left\{ \frac{\gamma(d - \tau)\alpha}{1 + \gamma}, \frac{d\alpha}{\gamma(1 + d\alpha) + 1} \right\}$$

where τ denotes the optimal value of threshold in a system with only fixed costs $F(\cdot)$.

Note that the bound in Theorem I.17 converges to 0 if either $\gamma \rightarrow 0$ or $\gamma \rightarrow \infty$ (the bound may not be tight for intermediate values). This is quite intuitive: when $\gamma \approx 0$, fixed costs dominate variable costs and the optimal policy is a threshold policy; when $\gamma \approx \infty$, variable costs dominate fixed costs and the optimal policy is to ship immediately upon order arrival (in other words, $\tau = d$).

1.7 Two Warehouses with Fixed and Variable Costs

In this section, we consider the general and practical setting introduced in Section 1.3, where the retailer delivers orders from at most two warehouses/stores and incurs both fixed and variable costs. Given the already complex structure of the optimal policy for the two-warehouse setting with only fixed cost (Section 1.5) and one-warehouse setting with both fixed and variable costs (Section 1.6), clearly the structure of the

³⁷This is based on the cost function fitted from UPS shipping rates; see Appendix A for more details. $\gamma = 0.2$ and α varies within (0.01,0.3).

optimal policy for two-warehouse setting with both fixed and variable costs will not be tractable. Therefore, inspired by results for the simplified cases discussed in previous sections, we first propose two base heuristics, warehouse-based and order-based heuristics, which are similar to the heuristics used in Section 1.5. Based on their performance, we modify the warehouse-based heuristic and label the modified one as the warehouse-based⁺⁺ (WB⁺⁺) heuristic. WB⁺⁺ performs close to optimal (with average optimality gap of 0.46%) and is easy-to-implement.

1.7.1 Base Heuristics

We discuss our base heuristics: *Warehouse-Based One-Package Heuristic* (WB) and *Order-Based One-Package Heuristic* (OB). The first heuristic is the same as the warehouse-based heuristic in Section 1.5. The second heuristic is similar to the order-based heuristic in Section 1.5, with a minor natural modification: In Section 1.5, since only fixed cost is considered, consolidation is decided by the comparison between $F_1(\min\{z_A, z_B\}) + F_2(z_C)$ and $F_1(z_A) + F_2(\min\{z_B, z_C\})$. To incorporate variable costs, consolidation is now decided by the comparison between $C_1(\min\{z_{A,1}, z_{B,1}\}, n_A + n_B) + C_2(z_{C,1}, n_C)$ and $C_1(z_{A,1}, n_A) + C_2(\min\{z_{B,1}, z_{C,1}\}, n_B + n_C)$, where $z_{X,1}$ denotes the smallest slack time for order type X and $C_i(z, n) = F_i(z) + n \cdot v_i(z)$. Details of simulations can be found in Appendix ??.

1.7.2 Performance of Base Heuristics: Warehouse-Based and Order-Based

In this section, we first analyze the performances of both warehouse-based and order-based heuristics using simulation and identify two reasons that could make them underperform. Based on these reasons, we propose an improved heuristic and numerically show the improvements.

We test the performances of WB and OB for two warehouses setting in a large scale

numerical experiments with 4,374 problem instances. We use the same set of parameters of customer-ordering frequency and fixed cost functions as in section 5.2. Variable costs are set as a portion of the fixed cost, with the same structure as commonly found in practice, $v_i(\cdot) = \gamma F_i(\cdot)$ ($i \in \{1, 2\}$). We vary $\gamma \in [0.1, 0.9]$ to cover a wide range of cases where the variable costs are relatively small compared to the fixed costs up to the cases where the variable costs are close to the fixed costs. The average simulation times to calculate the values of thresholds for the two heuristics are 0.12 and 0.06 seconds for $d = 3$, and 42.19 and 7.86 seconds for $d = 10$, respectively. We compare the performances of our proposed heuristics with the optimal cost (for some cases) and with three commonly used heuristics: (1) Myopic heuristic, which ships orders immediately upon arrival; (2) Time-threshold heuristic, which ships orders every several periods from each warehouse; and (3) Volume-threshold heuristic, which ships all pending orders of a certain type whenever the volume of that type triggers a threshold. While the policies used in practice may slightly vary, the dominating ones are simple myopic policies. We understand from Amazon.com employees and the Chinese online retailers that their current practices continue to follow the “old” rules, while considering alternatives to be implemented in the near future. Currently, the arrivals of orders are treated as a flow. Such flow is directed to the warehouse and each order is almost immediately packed and shipped. Thus, orders are shipped out (myopically) as early as operationally possible. We believe that the myopic shipment policy can be used as a benchmark. As for the time- and volume-threshold heuristics, as discussed in Section 1.2, they are the two most widely studied heuristics in the consolidation literature (Cooper, 1984; Higginson and Bookbinder, 1994, Cetinkaya et al., 2000, 2008). All three heuristics (myopic, time-threshold, and volume-threshold) are appropriately modified to allow consolidation across the warehouses and the parameters for each heuristic are optimized through simulation-based optimization.

We separately describe the results of our experiments for small and large size problems. For small size problems ($d \leq 5$), the optimal cost can be numerically computed and, thus, the costs of all heuristics are compared with the optimal cost. The average gaps across parameters other than d are reported in Table 1.2. Note that both warehouse-based and order-based heuristics are close to optimal and clearly outperform the other heuristics by significant margins.

d	Warehouse-based	Order-based	Myopic	Time-threshold	Volume-threshold
3	0.66%	0.11%	15.44%	15.44%	1.22%
4	1.60%	0.29%	20.16%	6.99%	3.32%
5	2.44%	1.43%	25.83%	6.58%	7.98%

Table 1.2: Percentage Gap to Optimal

For large size problems ($d > 5$), the optimal cost cannot be computed in a reasonable time. (Solving the DP for a single instance of $d = 5$ and $d = 6$ takes more than 5 hours and 20 days, respectively.) This motivates us to use different benchmarks to assess the performance of the heuristics. Given that the order-based heuristic performs close to optimal for $d \leq 5$, we measure the performances of all heuristics by the percentage gaps to the cost of order-based heuristic, see Table 1.3.³⁸

d	Warehouse-based	Myopic	Time-threshold	Volume-threshold
6	1.43%	37.59%	3.57%	7.41%
8	1.80%	60.22%	3.52%	8.52%
10	2.42%	80.79%	3.77%	9.52%

Table 1.3: Percentage Gap to Order-Based Heuristic

We find that there are two reasons that make the WB heuristic performs slightly worse than OB. These reasons provide directions for improvement:

(1) The warehouse-based consolidation rule effectively assigns $\max\{\tau_1, \tau_2\}$ as the threshold for orders of type B (as type B orders are shipped when either τ_1 or τ_2 is triggered). Since the value of this threshold is higher for the more expensive warehouse,

³⁸The reasons that the performance gaps of warehouse-based and order-based heuristics are functions of d are subtle. It may be interpreted as due to the length of due dates or due to the discretization of shipping options. We see that the discretization plays a big role: for each d , the cost function is discretized into d segments. With fewer choices (small d), the likelihood of choosing wrong threshold is lower.

orders type B cross the higher threshold first and are shipped from the more expensive warehouse. Consequentially, we consider the following modification: the threshold for orders type B is set to $\min\{\tau_1, \tau_2\}$, which allows orders type B to be shipped from the cheaper warehouse. In the numerical tests for the case $d = 4$, this modification decreases the optimality gap by 0.27% on average, from 1.60% to 1.33%).

(2) The warehouse-based consolidation rule ignores perfectly foreseeable opportunities to combine orders that are already placed and orders that will arrive in the future. Specifically, by looking at the slack times of orders of other types, one can intentionally delay some orders and ship them from the cheaper warehouse later.³⁹ Motivated by this observation, we replace the warehouse-based consolidation rule by a simple slack-time-dependent rule: the one-period myopic rule. In the numerical tests for the case $d = 4$, this improvement decreases the optimality gap by 0.94% on average, from 1.33% to 0.39%.

In addition to considering improvements for WB heuristic, we also consider potential improvement to OB heuristic. We try other slack-time dependent consolidation rules beyond the one-period rule. However, in the numerical tests, we observe that the one-period rule performs better than or almost as good as other rules we considered.⁴⁰ Thus, we continue using the one-period myopic consolidation rule in OB heuristic.

³⁹Consider the case where $d = 3$, $F_1 = F_2$, $v_1 = 0.9 * F_1$, $v_2 = 0.1 * F_2$, $\alpha = (0.5, 0.5, 0.5)$, and the optimal warehouse-based heuristics are $(\tau_1, \tau_2) = (3, 2)$. For state $(z_A, z_B, z_C) = (3, 3, 3)$, orders of type A trigger the thresholds τ_1 , by the warehouse consolidation rule, orders of type B is shipped from warehouse 1. However, it is easy to see that shipping B with C from warehouse 2 in the next period actually incurs less cost.

⁴⁰We first consider both the current one-period myopic consolidation rule and a two-period rule, and compare their performances. Observing that neither one dominates the other, we consider separating rules to utilize both one-period and two-period rule in the heuristic. However, even for the best separating rules we tested, it over performs the one-period rule by only 0.05% on average.

1.7.3 Warehouse-based⁺⁺ Heuristic

In this section, we modify the WB heuristic to the warehouse-based⁺⁺ heuristic (WB⁺⁺). Based on observations in the previous subsection, WB⁺⁺ uses thresholds τ_1 and τ_2 , which are similar to the WB. For better performance, WB⁺⁺ uses $\min\{\tau_1, \tau_2\}$ for orders type B and also uses the myopic consolidation rule. However, these features still mean that it is hard to directly solve the optimal value of τ 's. Thus, we introduce a direct way to calculate the value of thresholds, which makes the heuristic very easy-to-implement.⁴¹

WB⁺⁺ Heuristic

Given (τ_1, τ_2) and $(z_{A,1}, z_{B,1}, z_{C,1})$:

1. Ship all orders of types X if: $z_{A,1} \leq \tau_1$ for $X = A$, $z_{B,1} \leq \min\{\tau_1, \tau_2\}$ for $X = B$, or $z_{C,1} \leq \tau_2$ for $X = C$.
2. Consolidate with orders of other types depending on the one-period myopic rule.

The direct way to calculate the value of (τ_1, τ_2) is as follows: (1) we assume that a certain portion (ω) of arrivals of B is directed to warehouse 1 and the remaining portion ($1 - \omega$) to warehouse 2. (2) We optimize the partition (ω) and the thresholds (τ_1, τ_2) to minimize the total costs of two separate warehouses.⁴² The formulation of the optimization is as below:

$$\min_{\omega} \left[\min_{\tau_1} [C_1(\alpha_A, \alpha_B \omega, \tau_1)] + \min_{\tau_2} [C_2(\alpha_B, \alpha_C(1 - \omega), \tau_2)] \right] \quad (1.4)$$

⁴¹Although we do not anticipate the need to apply it for more than two warehouses (the practical setting support that using two warehouses will likely be the rule for the foreseeable future), the WB⁺⁺ heuristic is easily scalable to multiple warehouses.

⁴²The benefit of using ω is to translate a complicated problem to two very simple separate-warehouse problems. (While the parameters are found by exhaustive search, this is computationally very efficient as single-warehouse problems are very easy to solve.) After identifying ω that results in the lowest cost, we only use the corresponding τ_1 and τ_2 , and do not use ω any further.

where $C_i(\alpha, \beta, \tau)$ is the average cost per period for a separate warehouse i ($i \in 1, 2$). Note that $C_i(\alpha, \beta, \tau) = \frac{F(\tau) + \mathbb{E}[n] * \nu(\tau)}{\mathbb{E}[L]}$ where $\mathbb{E}[L] = \frac{1}{2\alpha\beta + (1-\alpha)\beta + (1-\beta)\alpha} + d - \tau$ is the expected number of periods between two consecutive shipments,⁴³ $\mathbb{E}[n] = \frac{2\alpha\beta}{M} + \frac{(1-\alpha)\beta}{M} + \frac{(1-\beta)\alpha}{M} + (d - \tau)(2\alpha\beta + (1 - \alpha)\beta + (1 - \beta)\alpha)$ is the average number of orders in a cycle, and $M = \alpha\beta + (1 - \alpha)\beta + (1 - \beta)\alpha$.⁴⁴

The intuition of the threshold calculation is as follows: the thresholds in WB⁺⁺ effectively assign a fraction of orders of type B to be shipped from warehouse 1 and the remaining fraction to be shipped from warehouse 2, depending on the slack times in each period. In the direct calculation, we replace the slack-time-based assignment with a constant assignment for the computational simplicity. Note that it still captures the economics of the cost structure and order arrival probabilities in each of the two-separate warehouses. Numerically tests⁴⁵ also show that the thresholds calculated in WB⁺⁺ perform very close to the optimized thresholds.⁴⁶

d	WB ⁺⁺	d	WB ⁺⁺
3	0.17%	6	0.09%
4	0.37%	8	0.20%
5	1.83%	10	0.81%

Table 1.4: Percentage Gap to Optimal

Table 1.5: Percentage Gap to Order-Based Heuristic

Comparing with the optimal policy, the WB⁺⁺ heuristic with directly-calculated thresholds also has very good performance. The performance of WB⁺⁺ is tested using the same set of parameters as in Section 1.7.2 and the results are summarized in Tables 1.4 and 1.5. In most of the cases, WB⁺⁺ performs well with optimality gap at most

⁴³Note that $\frac{1}{2\alpha\beta + (1-\alpha)\beta + (1-\beta)\alpha}$ is the number of periods until the first order arrived and $d - \tau$ is the number of periods until the threshold is triggered.

⁴⁴ $\frac{2\alpha\beta}{M} + \frac{(1-\alpha)\beta}{M} + \frac{(1-\beta)\alpha}{M}$ is the expected number of orders in the period, when the first order arrives (there are two streams of orders; α and β), and $(d - \tau)(2\alpha\beta + (1 - \alpha)\beta + (1 - \beta)\alpha)$ is the average number of orders in the following $(d - \tau)$ periods.

⁴⁵We tested on the additive cases where $F_1 = F_2 + \beta_1$. The multiplicative case has a similar result.

⁴⁶The average percentage optimality gap of the optimized τ 's is 0.38% (Stdv: 0.43%, Max: 3.67%), while the average percentage optimality gap of the solved τ 's in WB⁺⁺ is 0.50% (Stdv: 0.64%, Max: 6.46%).

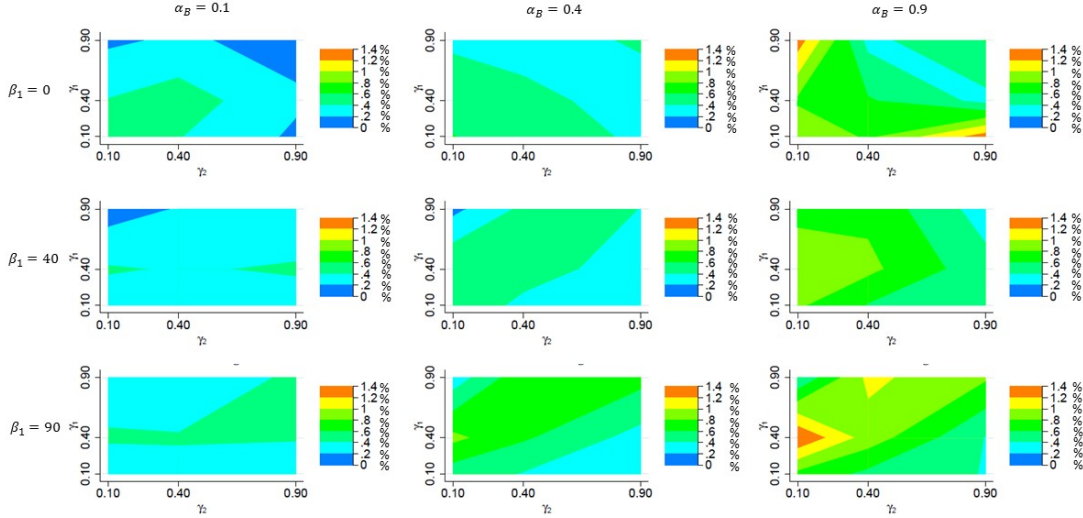


Figure 1.5: Average percentage gaps across different d , α_A , and α_C

2% (the cases with gaps larger than 2% accounts for about 4% of the total cases tested). To better understand its performance, we also visualize the optimality gaps of the cases with $F_2 = F_1 + \beta_1$ ⁴⁷ and $v_i = \gamma_i F_i$ $i \in \{1, 2\}$, in Figure 1.5 and 1.6. We identify two conditions where the heuristic can perform poorly: The first condition, as is shown in Figure 1.5, is where the variable costs dominate in one warehouse while the fixed costs dominate in the other warehouse, with a high arrival probability of orders type B. Specifically, this corresponds to the cases where one warehouse has relatively larger variable costs but a relatively smaller or equal fixed costs, compared to the other warehouse. Some examples of such cases would be: (1) warehouse 1 has large variable costs, $\gamma_1 = 0.9$, and warehouse 2 has small variable costs, $\gamma_2 = 0.1$, but both warehouses have the same fixed costs; (2) the opposite case, with $\gamma_1 = 0.1$ and $\gamma_2 = 0.9$; and (3) the case where warehouse 2 has larger fixed costs ($\beta_1 = 90$) and warehouse 1 has larger variable costs ($\gamma_1 = 0.4$ compared with $\gamma_2 = 0.1$). The second condition, as is shown

⁴⁷The cases with $F_2 = \beta_2 F_1$ has similar results.

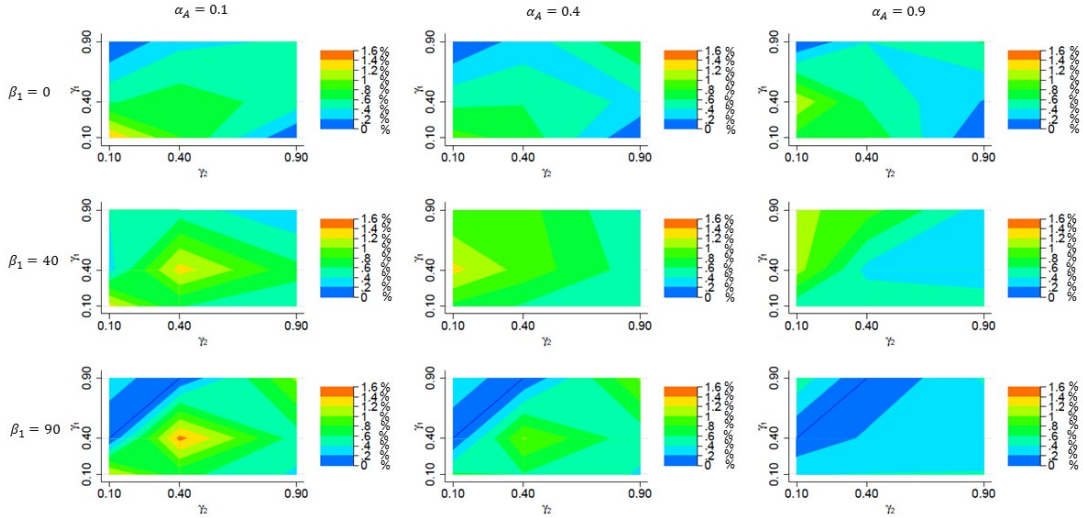


Figure 1.6: Average percentage gaps across different d and α_C , with $\alpha_B = 0.9$

in Figure 1.6,⁴⁸ is where both the fixed and variable costs of one warehouse are either similar to or lower than that of the other warehouse, but the product that can be only fulfilled from this warehouse has a fairly small arrival probability. An example would be the case where warehouse 1 has both smaller fixed costs ($\beta_1 = 90$) and smaller variable costs ($\gamma_1 = 0.1$ and $\gamma_2 = 0.9$), but product A has an arrival probability of only 0.1. Note that these two conditions are necessary but not sufficient.

The poor performance of the heuristic is driven by a mismatch between the pre-determined allocation in the calculation of thresholds and the state-dependent allocation when utilizing the thresholds. Such mismatch is exaggerated in the two conditions identified above. For the first condition, the pre-determined allocation rule tends to

⁴⁸As is shown below, the two conditions we identified are mutually exclusive. So we can first exclude the cases which are already identified in the first condition (and with gaps more than 1.2%), then plot the remaining data in this graph. Also, note that as we plot data from a different “angle” (taking average across different parameters) in Figure 1.6, it contains cases that met the first condition but did not appear to perform poorly in Graph 1.5, e.g., the case with $\beta_1 = 40$, $\gamma_1 = 0.4$, $\gamma_2 = 0.1$, and $\alpha_A = 0.4$. Further, as a robustness check, we also plot the data from a third “angle” after deleting the cases that met the first two conditions – the poor performance area does not show up anymore.

assign a higher-than-necessary threshold for the warehouse with larger variable costs while assigning a lower-than-necessary threshold for the warehouse with larger fixed costs. This is because the thresholds are calculated as if they are for two separate warehouses—less orders of type B are pre-allocated to the warehouse with higher variable costs, which pushes the calculated thresholds higher.⁴⁹ For the second condition, the pre-allocation rule allocates all orders of type B to the cheaper warehouse and thus assigns a lower-than-necessary threshold for the more expensive warehouse (when a warehouse does not have dominating variable or fixed costs, it is easy to observe from the formulation of $C_i(\alpha, \beta, \tau)$ that the optimal threshold increases with β), while in the actual slack-time-dependent allocation, there are, of course, orders that need to be shipped from the more expensive warehouse.⁵⁰

In reality, however, the parameters that meet the two conditions above rarely occur: (1) warehouses with higher fixed costs usually also have larger variable costs,⁵¹ and (2) the cheaper warehouses are usually made to fulfill more orders, which means higher arrival probability.

It also worth noting that there are cases where the retailer can simply use the naive heuristics and achieve good performance as well as the WB^{++} heuristic: when both warehouses have very high fixed costs (variable costs), then obviously it is optimal to hold orders as late as possible (to ship out orders as early as possible), where the WB^{++} heuristic is effectively equivalent to the case with $(\tau_1, \tau_2) = (1, 1)$ (the myopic policy with $(\tau_1, \tau_2) = (d, d)$). In such cases, naive heuristics perform well. However, these cases are extreme and are not the majority cases in practice. In other cases, the naive benchmark heuristics either assign wrong values of thresholds (e.g., myopic

⁴⁹For example, in the case of $\beta_1 = 90$, $\gamma_1 = 0.4$, $\gamma_2 = 0.1$ and $(\alpha_A, \alpha_B, \alpha_C) = (0.1, 0.9, 0.1)$, the optimal thresholds are (3, 2) while the thresholds in WB^{++} are (4, 1).

⁵⁰For example, in the case of $\beta_1 = 90$, $\gamma_1 = 0.4$, $\gamma_2 = 0.4$ and $(\alpha_A, \alpha_B, \alpha_C) = (0.1, 0.9, 0.4)$, the optimal thresholds are (3, 2) while the thresholds in WB^{++} are (3, 1).

⁵¹Many shipping costs are distance-based, e.g. shipping via UPS from nearby warehouses incur both lower fixed cost and variable cost per pound as opposed to a warehouse farther away.

policy), or use the wrong triggers which mess up with the slack-based delivery schedule (e.g., time-based and volume-based policy.)

Overall, while OB performs slightly better than WB^{++} , WB^{++} is easier to implement in terms of both parameter estimation and execution—the value of the thresholds can be calculated directly from optimization (1.4) and the warehouse-wise thresholds are easy to execute. In practice, both OB and WB^{++} are effective options for retailers, who can choose which heuristic to use.

To implement these heuristics in practice, the parameters need to be adequately chosen. (1) The length of period should be set to reflect the available shipping modes. For example, in omni-channel retail, where the shipping modes are 20 mins, 40 mins, 60 mins, etc, the time unit would naturally be set as 20 mins. In online retail, where the shipping modes are same-day, next-day, three-day, etc, the time unit would be set to a day. (2) The order arrival probabilities ($\alpha_{X,m}$, $X \in \{A, B, C\}$) should be calculated based on relevant historical data. For omni-channel, the arrival probability can be set as the average order frequency of customers from the same region. For online retailer, since each customer is treated separately, the arrival probability can be set according to his average frequency of placing orders. (3) The horizon length is not needed for heuristic calculations.

1.8 Conclusion

This paper analyzes emerging and increasingly promising areas of consolidation of orders in outbound logistics. In addition to its general applicability in the supply chain setting, it can also be a potentially efficient way to improve the current practice of outbound shipment in Business-to-Customer settings, including e-commerce and omni-channel retail. Consolidating orders placed at different times can reduce the number of shipments and decrease the total shipping costs. However, it leads to new trade-offs

that retailers need to carefully balance: on one hand, combining several orders into one shipment can eliminate some of the fixed shipping costs; on the other hand, it may cause a shipping delay to consolidate current orders with future one and, thus, may require expedited more-expensive shipping. In this paper, we have analyzed the optimal consolidation policy and described its structure for the practical setting with two warehouses. We focus on crucial factors that determine the economies of the situations: fixed cost, variable cost, options to expedite, and different product availabilities. The structure of the optimal policy in three simplified cases is derived: (1) For one warehouse with fixed cost, we show that the optimal policy can be directly characterized by a sequence of time-dependent thresholds. (2) For two warehouses, with fixed cost and overlapping availability of products, the optimal policy, in the symmetric settings, can be characterized by six non-linear boundaries in three-dimensional space. We show that in asymmetric settings, the optimal policy can be approximated by constant-threshold heuristics, with less than a 2% optimality gap. (3) For one warehouse and both fixed and variable costs, the optimal policy is a function of volume-dependent thresholds, which can be approximated by constant thresholds, with less than a 0.3% optimality gap. Based on these special cases, three easily implementable threshold-based heuristic policies, WB, OB, and WB^{++} , are proposed, which significantly outperform the best benchmark policies found in practice/literature. We show that among them, WB^{++} and OB can be very useful in practice. The OB is very efficient, with 0.61% optimality gap on average, and the WB^{++} is easy to implement and efficient, with 0.79% optimality gap on average.

CHAPTER II

On a Deterministic Approximation of Inventory Systems with Sequential Probabilistic Service Level Constraints

2.1 Introduction

The problem of minimizing inventory cost over time while providing a high quality customer service in the presence of stochastic demand is one of the most fundamental and challenging core problems of inventory management. Depending on whether unmet demand can be satisfied at a later time or not, many inventory systems can be categorized as either a backorder or lost-sales system. In the so-called canonical *cost-based* model where there is no service level constraint (i.e., no explicit targeted service level guarantee) and the objective is simply to minimize the expected total ordering, holding, and stock-out costs, both backorder and lost-sales inventory models have been extensively studied in the literature (cf. Zipkin [96]). It is a common belief that, in the presence of positive lead time (i.e., the time lag between when an order for more inventory is placed and when the order is received), it is much more difficult to optimize a lost-sales inventory system than its backorder counterpart. Indeed, while the backorder system usually has a simple optimal order-up-to (or base-stock) type of control

(e.g., Zipkin [96]), the lost-sales system is notoriously challenging to analyze and the structure of its optimal control is poorly understood. We refer interested reader to an excellent review paper by Bijvank and Vis [15] for more discussions on lost-sales inventory models. As an exact solution seems out of reach, asymptotic analysis has been performed recently by various researchers. For example, Huh et al. [47] show that an order-up-to control is asymptotically optimal for the lost-sales system as the lost-sales penalty grows. The intuition of this result is that the event of stock-out becomes rare as the lost-sales penalty grows; in fact, so much so that the lost-sales system behaves nearly identical to its corresponding backorder system. Hence, the order-up-to control that is optimal for the backorder system is asymptotically optimal for the lost-sales system. Note that although this argument may seem quite intuitive, it does not directly translate into a simple mathematical proof. To the contrary, the proof of this seemingly intuitive result is non-trivial. On a different asymptotic regime, a simple constant-order policy that places an order with exactly the same quantity in every time period, regardless of the current inventory level, is proved to be asymptotically optimal for the lost-sales system as the lead time grows in Goldberg et al. [40], Xin and Goldberg [89]. The intuition of this result is that, since lead time is large, there is a significant amount of randomness in the system between when an order is placed and when the order is eventually received. This suggests that dynamically adjusting order quantities over time may not provide much benefit compared to a passive control that simply order the same quantity at every time period. Perhaps this argument seems intuitive as well, it does not directly translate into a simple mathematical proof. Other recent works in the inventory literature that use asymptotic analysis include Reiman and Wang [67] (large lead time asymptotic), Wan and Wang [86] (large volume asymptotic), Ahn et al. [3] (large batch size asymptotic), Bu, Gong and Yao [22] (large lead time asymptotic), and Xin [91] (large lead time asymptotic).

While much of the inventory literature in the past decade have focused on the cost-based model for tractability reason, there are also works in the literature that study inventory control under a so-called *service-based* model. This is motivated by the fact that, in practice, the cost of unmet demand is often difficult to quantify (e.g., Chen and Krass [30]) and, therefore, service level is typically used as a more direct metric for evaluating the performance of inventory replenishment controls (see Bertsimas and Paschalidis [13] for more discussions on the drawback of cost-based model). For example, some firms such as Walmart use both average on-hand inventory and service level as their two key performance metrics (Xin et al. [90]). In the literature, there is more than one way to define a service level. The event-oriented α -service level is defined as the probability of no stock-out; the quantity-based β -service level is defined as the proportion of total demand that is immediately satisfied without delay, capturing not only the stock-out event but also the amount of stock-out; the time-and-quantity-related γ -service level is defined to reflect not only the amount of backorders but also the waiting times of the demands backordered. Among these three, the α -service level is one of the most widely used service level criteria in the inventory literature (e.g., Snyder and Shen [71]) and is also recognized in practice (see Jiang, Shi and Shen [50] for practical examples). In our paper, we focus on the sequential version of the α -service level. Note that if the service-level constraint is not sequential (as considered in Bitran and Yanasse [19]), such a constraint is essentially an expected service-level requirement which is set upfront. However, in practice, when a company makes decisions in week t , all realized demand information in the past are naturally considered to meet the in-stock probability requirement in the following weeks. The sequential service-level constraints we considered here appropriately reflect such a practice.

Our work is closely related to some of the early studies of periodic review inventory models with α -service level constraints such as Bitran and Yanasse [19], Bitran and

Sarkar [18], Bitran and Leong [17]. The authors in these papers assume a backorder inventory system and use deterministic programs to approximate the stochastic inventory problems. They derive bounds on the gap between the two formulations to show that deterministic inventory models well-approximate stochastic inventory models in the regime of high service level (i.e., the setting where probability of stock-out is close to zero). Their results have an important practical implication: while the multi-period stochastic inventory problem with service level constraints is very difficult to solve, decision-makers can use its deterministic approximation for the purpose of estimating total costs in the context of strategic inventory planning. In this paper, we ask whether similar results can be established for a more complex lost-sales inventory system. Specifically, we consider both the backorder and lost-sales inventory systems with positive lead time and sequential α -service level constraints, and analyze the performance of a simple order-up-to control in the regime of high service level requirement. The parameters of our heuristic control can be computed using the optimal solution of a deterministic program (in fact, a linear program), which is a deterministic approximation of the stochastic backorder system. We show that this order-up-to control is asymptotically optimal for both the backorder and lost-sales inventory systems as the service level increases to 100%. In summary, our results have two main contributions. First, it complements the results of Bitran and Yanasse [19], Bitran and Sarkar [18], and Bitran and Leong [17] by providing an asymptotically optimal order-up-to control for the backorder system. Specifically, while Bitran and Yanasse [19] and Bitran and Leong [17] consider the expected service-level constraint, we consider the more realistic sequential service-level constraint. While Bitran and Sarkar [18] use a deterministic system to approximate the sequential problem, it does not provide an asymptotically optimal heuristic as we do – this is mainly because their deterministic system does not incorporate the updating nature of the sequential constraints. Second, our work fur-

ther shows that the order-up-to control that we propose is also asymptotically optimal for the lost-sales system. In terms of methodological contributions, our analysis for the lost-sales system involves a construction of an alternative backorder system whose expected total cost can be related to that of the analogous lost-sales system. In terms of practical implications, our results give credence to the use of deterministic program to approximate complex lost-sales inventory problem with service level constraints. In fact, since many real-world inventory problems with service level constraints are difficult to solve, the majority of existing research has simply focused on directly analyzing the deterministic formulation of the problems (e.g., Tarim and Kingsman [75], Tarim et al. [76], Worm et al. [88]). Although our result is specific to the setting that we consider in this paper, we believe that it is an important step for further analyzing the quality of deterministic approximation in other related inventory systems with service level constraints (e.g., joint inventory and fulfillment decisions, and dual-sourcing problems).

It is worth noting that our work shares the same spirit as Huh et al. [47] in the sense that both papers show that the lost-sales model is asymptotically identical to its backorder counterpart as the lost-sales penalty (or, equivalently, the service level) increases. However, there are some notable differences. In Huh et al. [47], demands are assumed to be independent and identically distributed (i.i.d.) and there is no explicit rate of convergence; by contrast, in our work, demands across different periods can be highly non-stationary, as long as they share the same support (see the definition in Section 2), and we also derive an explicit bound for the optimality gap. In Huh et al. [47], the order-up-to level is derived from the stochastic backorder system whereas the order-up-to level in our heuristic control is computed using the optimal solution of a linear program. Overall, our results in the service-based model complement their results in the cost-based model.

Aside from the above cited papers that use deterministic programs to approximate the original stochastic inventory problem, there are other works in the literature that study inventory problem with service level constraints with the focus on developing/analyzing heuristic controls of certain forms. These include Boyaci and Gallego [21], Shang and Song [68], Özer and Xiong [63], Bijvank and Vis [16], and Bijvank [14]. More recently, Jiang, Shi and Shen [50] study an inventory model with backorder and α -service level constraints, and propose a novel 2-approximation algorithm based on the idea of delayed forced holding and production cost.

2.1.1 Outline of Paper

The rest of the paper is organized as follows. We formally define our inventory models with service level constraints and lead times and introduce a deterministic backorder relaxation in Section 2.2. We propose our heuristic controls, state our main results, and discuss the key ideas of the proofs for the backorder and lost-sales systems in Section 2.3 and Section 2.4 respectively (unless otherwise noted, all remaining details of the proofs can be found in the e-companion of the paper). We conclude and propose directions for future research in Section 2.5. We also attach a technical appendix in the e-companion.

2.2 Model Description

We consider a multi-period stochastic inventory problem with independent but possibly non-stationary (time-varying) demands. Demand at period t is denoted by D_t and has a support in interval $[0, \bar{D}]$ for some $\bar{D} < \infty$. Specifically, we assume that $P(D_t \in [0, \bar{D}]) = 1$ and $P(D_t \leq \bar{D} - \epsilon) < 1$ for all $\epsilon > 0$. Lead time is deterministic and equals $L \geq 0$. We consider both backorder and lost-sales systems. As usual, we use c to denote the per-unit ordering cost, h to denote the per-unit per-period holding

cost, and p to denote either the per-unit per-period backorder cost or the per-unit per-period lost-sale penalty, depending on whether the corresponding system is backorder or lost-sales. (Although we study service-based model, we still include penalty cost to allow instances with real penalty cost (i.e., the cost of drop-shipping, etc.). That said, consistent to the philosophy of service-based model, in our asymptotic analysis, we will assume that p is fixed while the value of service level is increased to 1.)

Let $I_t^{\pi,b}$ denote the inventory level at the beginning of period t in the backorder system under control π before the new order arrives. Since lead time equals L , we have: $I_{t+1}^{\pi,b} = I_t^{\pi,b} + x_{t-L}^{\pi,b} - D_t$, where $x_t^{\pi,b}$ is the quantity ordered under control π in period t . Similarly, let $I_t^{\pi,\ell}$ denote the inventory level at the beginning of period t in the lost-sales system under control π before the new order arrives. Under the lost-sales system, we have: $I_{t+1}^{\pi,\ell} = (I_t^{\pi,\ell} + x_{t-L}^{\pi,\ell} - D_t)^+$ where $x_t^{\pi,\ell}$ is the quantity ordered under control π in period t . Let $\mathfrak{S}_t^{\pi,b}$ and $\mathfrak{S}_t^{\pi,\ell}$ denote the history of all realizations (both demands and ordering decisions) up to the end of period t under control π for the backorder and lost-sales systems, respectively. Also, let $C^{*,b}(\alpha)$ and $C^{*,\ell}(\alpha)$ denote the optimal expected total costs over a finite horizon under the backorder and lost-sales systems with service level $1 - \alpha$, respectively. We can write $C^{*,b}(\alpha)$ and $C^{*,\ell}(\alpha)$ as follows:

$$\begin{aligned}
C^{*,b}(\alpha) &= \min_{\pi \in \Pi^b} \mathbf{E} \left[\sum_{t=1}^T c \cdot x_{t-L}^{\pi,b} + \sum_{t=1}^T h \cdot (I_t^{\pi,b} + x_{t-L}^{\pi,b} - D_t)^+ \right. \\
&\quad \left. + \sum_{t=1}^T p \cdot (D_t - I_t^{\pi,b} - x_{t-L}^{\pi,b})^+ \right], \tag{2.1} \\
C^{*,\ell}(\alpha) &= \min_{\pi \in \Pi^\ell} \mathbf{E} \left[\sum_{t=1}^T c \cdot x_{t-L}^{\pi,\ell} + \sum_{t=1}^T h \cdot (I_t^{\pi,\ell} + x_{t-L}^{\pi,\ell} - D_t)^+ \right. \\
&\quad \left. + \sum_{t=1}^T p \cdot (D_t - I_t^{\pi,\ell} - x_{t-L}^{\pi,\ell})^+ \right],
\end{aligned}$$

where the set of feasible controls Π^b and Π^ℓ are defined as:

$$\begin{aligned}\Pi^b &= \left\{ \pi : x_t^{\pi,b} \geq 0, P \left(I_t^{\pi,b} + x_{t-L}^{\pi,b} - D_t \geq 0 \right) \geq 1 - \alpha, \quad \forall 1 \leq t \leq L + 1 \right. \\ &\quad \left. \text{and } P \left(I_t^{\pi,b} + x_{t-L}^{\pi,b} - D_t \geq 0 \mid \mathfrak{S}_{t-L-1}^{\pi,b} \right) \geq 1 - \alpha, \quad \forall L + 2 \leq t \leq T \right\}, \\ \Pi^\ell &= \left\{ \pi : x_t^{\pi,\ell} \geq 0, P \left(I_t^{\pi,\ell} + x_{t-L}^{\pi,\ell} - D_t \geq 0 \right) \geq 1 - \alpha, \quad \forall 1 \leq t \leq L + 1 \right. \\ &\quad \left. \text{and } P \left(I_t^{\pi,\ell} + x_{t-L}^{\pi,\ell} - D_t \geq 0 \mid \mathfrak{S}_{t-L-1}^{\pi,\ell} \right) \geq 1 - \alpha, \quad \forall L + 2 \leq t \leq T \right\}.\end{aligned}$$

Both $C^{*,b}(\alpha)$ and $C^{*,\ell}(\alpha)$ depend on T and we suppress this dependency on our notations. Note that if $p = 0$, both backorder and lost-sales problems become trivial and an optimal control orders as little as possible as long as the service level constraint is satisfied.

Assumptions. We make the following modeling assumptions:

- A1. The order quantities arriving in periods $1, 2, \dots, L, L+1$ (i.e., $x_{1-L}, x_{2-L}, \dots, x_0, x_1$) are decided jointly at the beginning of period 1;
- A2. The initial inventory level at the beginning of period 1 is zero, i.e., $I_1^{\pi,b} = I_1^{\pi,\ell} = 0$ for all π ;
- A3. We allow fractional (continuous) order quantities and demand fulfillment.

Assumptions A1 and A2 are made without loss of generality and are useful to simplify some of the analysis. For assumption A1, we can alternatively assume that the quantities arriving in periods 1 to $L+1$ have been decided beforehand. In this case, since we have no control over these quantities, the sequence of probabilistic service level constraints defined in Π^b and Π^ℓ will have to start from period $L+2$ instead of period 1. Assumption A3 is a standard assumption made in the inventory literature and is often made to simplify the analysis. In addition, the impact of rounding error due to

integrality consideration is relatively negligible when order quantities are on the scale of hundreds or more.

On cost-based vs. service-based model. One advantage of our service-based model is that the control derived from our model can always achieve the required service level. However, this is not true for the canonical cost-based model as its objective is to simply minimize the expected total costs. We demonstrate this point through the following example when demands are stochastically decreasing by showing that the service level under an optimal solution of the canonical cost-based model can be arbitrarily poor. Suppose that $c = 0$, $L = 0$, $T = 2$, $h = 1$, $p = \frac{1-\alpha}{\alpha}$ (to ensure $\frac{p}{h+p} = 1 - \alpha$) and $D_2 = 0$ w.p.1. Then, the optimal order-up-to levels of the cost-based model are $x_2^c = 0$ and

$$x_1^c \in \arg \min_{x_1 \in \mathbb{R}} \mathbb{E} \left[h \cdot (x_1 - D_1)^+ + p \cdot (D_1 - x_1)^+ + h \cdot (x_1 - D_1)^+ \right].$$

The above is equivalent to $x_1^c = F_{D_1}^{-1} \left(\frac{p}{p+2h} \right) = F_{D_1}^{-1} \left(\frac{1-\alpha}{1+\alpha} \right)$. It can only achieve service level $\frac{1-\alpha}{1+\alpha}$ in the first period, which is less than $1 - \alpha$. Similarly, it is not difficult to extend it to general T with $D_t = 0$ w.p.1. for $t = 2, \dots, T$. In that case, the service level in the first period is $\frac{1-\alpha}{1+(T-1)\alpha}$, which can be arbitrarily small for a sufficiently large T .

A Deterministic Approximation of Backorder System. We now discuss a deterministic approximation of backorder system. The solution of this deterministic formulation will be later used to construct heuristic controls for the original (stochastic) backorder and lost-sales systems.

First, note that, in the backorder system, we can write:

$$I_t^{\pi,b} + x_{t-L}^{\pi,b} - D_t = \begin{cases} I_1^{\pi,b} + \sum_{s=1}^t x_{s-L}^{\pi,b} - \sum_{s=1}^t D_s \quad \forall 1 \leq t \leq L+1; \\ I_{t-L}^{\pi,b} + \sum_{s=t-L}^t x_{s-L}^{\pi,b} - \sum_{s=t-L}^t D_s \quad \forall L+2 \leq t \leq T. \end{cases} \quad (2.2)$$

Let $\mu_t = \mathbf{E}[D_t]$ denote the expected demand in period t , and define $\beta_t^k(\alpha)$ to be the smallest $\beta \in [-\sum_{s=t}^{t+k-1} \mu_s, \infty]$ satisfying inequality

$$P\left(\beta + \sum_{s=t}^{t+k-1} \mu_s - \sum_{s=t}^{t+k-1} D_s \geq 0\right) \geq 1 - \alpha$$

for $\alpha \in [0, 1]$. That is, $\sum_{s=t}^{t+k-1} \mu_s + \beta_t^k(\alpha)$ can be interpreted as the $(1 - \alpha)$ -quantile of $\sum_{s=t}^{t+k-1} D_s$. Using $\{\beta_t^k(\alpha)\}$, the probabilistic service level constraints in Π^b can be equivalently written as linear constraints:

$$I_1^{\pi,b} + \sum_{s=1}^t x_{s-L}^{\pi,b} - \sum_{s=1}^t \mu_s \geq \beta_1^t(\alpha) \quad \forall 1 \leq t \leq L+1, \quad (2.3)$$

$$I_{t-L}^{\pi,b} + \sum_{s=t-L}^t x_{s-L}^{\pi,b} - \sum_{s=t-L}^t \mu_s \geq \beta_{t-L}^{L+1}(\alpha) \quad \forall L+2 \leq t \leq T. \quad (2.4)$$

Constraints (2.3) and (2.4) motivate us to define the following deterministic model:

$$\begin{aligned}
\mathbf{DET:} \quad D^{*,b}(\alpha) = & \min_{x,y,z,m} \sum_{t=1}^T c \cdot x_{t-L} + \sum_{t=1}^T h \cdot z_{t+1} + \sum_{t=1}^T p \cdot m_{t+1} & (2.5) \\
\text{s.t.} \quad & y_1 = 0 \\
& y_t = y_{t-1} + x_{t-1-L} - \mu_{t-1} & \forall 2 \leq t \leq T+1 \\
& z_t \geq y_t & \forall 2 \leq t \leq T+1 \\
& m_t \geq -y_t & \forall 2 \leq t \leq T+1 \\
& \sum_{s=1}^t x_{s-L} - \sum_{s=1}^t \mu_s \geq \beta_1^t(\alpha) & \forall 1 \leq t \leq L+1 \\
& \sum_{s=1}^t x_{s-L} - \sum_{s=1}^t \mu_s \geq \beta_{t-L}^{L+1}(\alpha) & \forall L+2 \leq t \leq T \\
& x_t, z_t, m_t \geq 0 & \forall 1-L \leq t \leq T-L.
\end{aligned}$$

The variables x_t , y_t , z_t , m_t in the above deterministic formulation can be interpreted as the order quantity placed in period t , the inventory level at the beginning of period t , the amount of inventory overstocked in the end of period t , and the amount of inventory understocked in the end of period t , respectively, all in a deterministic backorder system. We use $x^{D,b}(\alpha) = (x_t^{D,b}(\alpha))$ and $y^{D,b}(\alpha) = (y_t^{D,b}(\alpha))$ to denote the optimal solution of (2.5) (note that the variables z_t and m_t are redundant). It is not difficult to show that $x_t^{D,b}(\alpha)$ can be written recursively as a function of $x_{t-1}^{D,b}(\alpha)$, $x_{t-2}^{D,b}(\alpha)$, \dots , $x_{1-L}^{D,b}(\alpha)$:

Lemma II.1. *We can write:*

$$\begin{aligned}
x_{1-L}^{D,b}(\alpha) &= \mu_1 + \beta_1^1(\alpha), \\
x_{t-L}^{D,b}(\alpha) &= \left(\sum_{s=1}^t \mu_s + \beta_1^t(\alpha) - \sum_{s=1}^{t-1} x_{s-L}^{D,b}(\alpha) \right)^+ & \forall 2 \leq t \leq L+1, \\
x_{t-L}^{D,b}(\alpha) &= \left(\sum_{s=1}^t \mu_s + \beta_{t-L}^{L+1}(\alpha) - \sum_{s=1}^{t-1} x_{s-L}^{D,b}(\alpha) \right)^+ & \forall L+2 \leq t \leq T.
\end{aligned}$$

It is not difficult to see from Lemma 1 that if demands are i.i.d with mean μ , then

$$\begin{aligned} x_{1-L}^{D,b}(\alpha) &= \mu + \beta_1^1(\alpha), \\ x_{t-L}^{D,b}(\alpha) &= \mu + \beta_1^t(\alpha) - \beta_1^{t-1}(\alpha), \quad \forall 2 \leq t \leq L+1, \\ x_{t-L}^{D,b}(\alpha) &= \mu, \quad \forall L+2 \leq t \leq T \end{aligned}$$

(for i.i.d. demands, $\beta_{t-L}^{L+1}(\alpha) = \beta_{t-L-1}^{L+1}(\alpha)$ for all $L+2 \leq t \leq T$). For the case with general independent demands, as long as α is sufficiently small, we have:

$$\begin{aligned} x_{1-L}^{D,b}(\alpha) &= \mu_t + \beta_1^1(\alpha), \\ x_{t-L}^{D,b}(\alpha) &= \mu_t + \beta_1^t(\alpha) - \beta_1^{t-1}(\alpha), \quad \forall 2 \leq t \leq L+1, \\ x_{t-L}^{D,b}(\alpha) &= \mu_t + \beta_{t-L}^{L+1}(\alpha) - \beta_{t-L-1}^{L+1}(\alpha), \quad \forall L+2 \leq t \leq T. \end{aligned}$$

Formulation **DET** is quite intuitive; indeed, as noted in Section 2.1, it has been widely used in the academic literature to approximate complex stochastic backorder systems with probabilistic service level constraints (e.g., Bitran and Yanasse [19], Bitran and Sarkar [18], Bitran and Leong [17]). The following lemma tells us that the optimal value of **DET** is a lower bound of the optimal expected total cost in the stochastic backorder system; thus, we can view **DET** as a relaxation of (2.1). This observation will be useful for our analysis later and we defer the proof to the appendix.

Lemma II.2. $D^{*,b}(\alpha) \leq C^{*,b}(\alpha)$.

2.3 Backorder Inventory System

In this section, we focus on backorder system. We first describe our heuristic control and discuss its theoretical performance in Section 3.1. Next, we provide an outline of

the proof of our main result (Theorem II.5) in Section 3.2. Finally, we report results from numerical experiments in Section 3.3.

2.3.1 Proposed Heuristic Control and Its Performance

We first describe our heuristic control for the backorder system, which we simply call H_b . Given the service level $1 - \alpha$ in period t , H_b places a new order $x_t^{H_b,b}(\alpha)$ defined as below:

$$x_t^{H_b,b}(\alpha) = \begin{cases} x_t^{D,b}(\alpha), & \forall 1 - L \leq t \leq 1; \\ \left(y_t^{D,b}(\alpha) + \sum_{s=t-L}^t x_s^{D,b}(\alpha) - I_t^{H_b,b} - \sum_{s=t-L}^{t-1} x_s^{H_b,b}(\alpha) \right)^+ & \\ \forall 2 \leq t \leq T - L. \end{cases} \quad (2.6)$$

Note that, by definition, H_b is an order-up-to control whose order-up-to level equals

$$y_t^{D,b}(\alpha) + \sum_{s=t-L}^t x_s^{D,b}(\alpha) \quad \forall t \geq 2.$$

Moreover, it is not difficult to see that H_b is a feasible control to the stochastic backorder system. (To see this, note that, from the definition of H_b , we have: $I_t^{H_b,b} + \sum_{s=t-L}^t x_s^{H_b,b} \geq y_t^{D,b}(\alpha) + \sum_{s=t-L}^t x_s^{D,b}(\alpha)$. Since $x^{D,b}(\alpha)$ and $y^{D,b}(\alpha)$ satisfy the deterministic analogue of constraints (2.3) and (2.4) in **DET**, consequently, constraints (2.3) and (2.4) are also satisfied by H_b . Hence, H_b satisfies the service level constraints and, therefore, it is feasible.) We state this formally as a lemma below.

Lemma II.3. $H_b \in \Pi^b$.

Since H_b is defined using the deterministic solution $x^{D,b}(\alpha)$, it is important to quantify the difference between the total orders placed by H_b in the stochastic system

and the total orders placed by the deterministic solution in the deterministic system as this difference will affect the performance of H_b (see the bound in Theorem II.5 below).

Lemma II.4. *We can bound:*

$$\mathbf{E} \left[\sum_{s=2}^{t-L} (x_s^{H_b,b}(\alpha) - x_s^{D,b}(\alpha)) \right] \leq \mathbf{E} \left[\sum_{s=1}^{t-L-1} (\mu_s - x_{s+1}^{D,b}(\alpha) - D_s)^+ \right] \quad \text{for all } t.$$

In particular, if demands are i.i.d, then

$$\mathbf{E} \left[\sum_{s=2}^{t-L} (x_s^{H_b,b}(\alpha) - x_s^{D,b}(\alpha)) \right] = 0 \quad \text{for all } t.$$

Observe that the first bound in Lemma II.4 goes to zero as $\alpha \rightarrow 0$. This is so because, per our discussions following Lemma II.1, for $t \geq 2$ and all sufficiently small α , we have:

$$\begin{aligned} x_t^{D,b}(\alpha) &= \mu_{L+t} + \beta_t^{L+1}(\alpha) - \beta_{t-1}^{L+1}(\alpha) \\ &= \mu_{t-1} + \left(\sum_{s=t}^{L+t} \mu_s + \beta_t^{L+1}(\alpha) \right) - \left(\sum_{s=t-1}^{L+t-1} \mu_s + \beta_{t-1}^{L+1}(\alpha) \right). \end{aligned}$$

Note that the two summations inside the (\cdot) converge to $(L+1)\bar{D}$ as $\alpha \rightarrow 0$. So, $x_t^{D,b}(\alpha) \rightarrow \mu_{t-1}$ as $\alpha \rightarrow 0$ for all $t \geq 2$ and, therefore, the first bound in Lemma II.4 goes to zero as $\alpha \rightarrow 0$. As for the second bound in Lemma II.4, it immediately follows from the first bound and our discussions following Lemma II.1 (i.e., the fact that $x_t^{D,b}(\alpha) = \mu$ for all $t \geq 2$ and all α).

The following theorem provides a bound for the loss of H_b with respect to the optimal control.

Theorem II.5. *We can bound:*

$$\begin{aligned}
C^{H_b,b}(\alpha) - C^{*,b}(\alpha) &\leq \sum_{t=2}^{T-L} [c + h \cdot (T - L - t + 1)] \cdot \mathbf{E} \left[x_t^{H_b,b}(\alpha) - x_t^{D,b}(\alpha) \right] \\
&\quad + (h + p) \cdot \sum_{t=1}^{L+1} \mathbf{E} \left[\left(\sum_{s=1}^t D_s - \sum_{s=1}^t \mu_s - \beta_1^t(\alpha) \right)^+ \right] \\
&\quad + (h + p) \cdot \sum_{t=L+2}^T \mathbf{E} \left[\left(\sum_{s=t-L}^t D_s - \sum_{s=t-L}^t \mu_s - \beta_{t-L}^{L+1}(\alpha) \right)^+ \right].
\end{aligned}$$

Moreover, we also have:

$$C^{*,b}(\alpha) \geq c \cdot \left[\sum_{t=1}^T \mu_t + \beta_{T-L}^{L+1}(\alpha) \right] + h \cdot \left[\sum_{t=1}^{L+1} \beta_1^t(\alpha) + \sum_{t=L+2}^T \beta_{t-L}^{L+1}(\alpha) \right].$$

Note that the term $\sum_{t=2}^{T-L} (T - L - t + 1) \cdot \mathbf{E}[x_t^{H_b,b}(\alpha) - x_t^{D,b}(\alpha)]$ in the bound can be expressed as $\sum_{k=2}^{T-L} \sum_{s=2}^k \mathbf{E}[x_s^{H_b,b}(\alpha) - x_s^{D,b}(\alpha)]$ and can, therefore, be further bounded using Lemma II.4. Specifically, if demands are i.i.d, the bound in Theorem II.5 can be bounded by a rough bound as follows:

$$C^{H_b,b}(\alpha) - C^{*,b}(\alpha) \leq (h + p) \cdot T \cdot (L + 1) \cdot \bar{D} \cdot \alpha,$$

which shows that $C^{H_b,b}(\alpha) - C^{*,b}(\alpha)$ converges to 0 (at least) linearly in α . If, on the other hand, demands are independent but not necessarily stationary, then in addition to the term $(h + p) \cdot T \cdot (L + 1) \cdot \bar{D} \cdot \alpha$, we also have the sum of $O(T^2)$ terms, each of which converges to 0 as $\alpha \rightarrow 0$.

2.3.2 Proof of Theorem II.5

We now provide an outline of the proof of Theorem II.5. For notational brevity, whenever there is no loss of information, and with an exception of a few notations such

as $\beta_t^k(\alpha)$, we will often suppress the dependency of all other notations on α (e.g., we will often write $x_t^{D,b}(\alpha)$ and $x_t^{H_b,b}(\alpha)$ simply as $x_t^{D,b}$ and $x_t^{H_b,b}$). Note that $C^{H_b,b} - C^{*,b}$ can be written as a sum of two terms:

$$C^{H_b,b} - C^{*,b} = [C^{H_b,b} - D^{*,b}] + [D^{*,b} - C^{*,b}].$$

By Lemma II.2, $D^{*,b} - C^{*,b} \leq 0$. To bound the second term, we proceed in several steps.

Step 1

Let $\Delta_t = D_t - \mu_t$ for all t and define

$$U_t^b := y_t^{D,b} - I_t^{H_b,b} + \sum_{s=t-L}^t (x_s^{D,b} - x_s^{H_b,b}) = \sum_{s=1-L}^t (x_s^{D,b} - x_s^{H_b,b}) + \sum_{s=1}^{t-1} \Delta_s$$

for $1 \leq t \leq T - L$. U_t^b represents the difference between the inventory position under the optimal policy of the deterministic model and that under policy H_b in period t . Also, let

$$k_t^b := \mathbf{E}[U_t^b - U_{t+1}^b] = \mathbf{E} \left[x_{t+1}^{H_b,b} - x_{t+1}^{D,b} \right]$$

for $1 \leq t \leq T - L - 1$, representing the difference between the order quantity of policy H_b and that of the optimal policy of the deterministic model in period $t + 1$.

The following lemmas are useful for our analysis (see Step 2 below).

Lemma II.6. *For $t \geq 2$, we have:*

$$x_t^{H_b,b} = (x_t^{D,b} + \Delta_{t-1} + U_{t-1}^b)^+ \quad \text{and} \quad U_t^b = -(U_{t-1}^b + \Delta_{t-1} + x_t^{D,b})^-$$

where $(x)^- = \max\{-x, 0\}$.

Lemma II.7. We can express $\mathbf{E}[I_t^{H_b,b}]$ as a function of $y_t^{D,b}$ and k_s ($s \leq t-1-L$) as follows:

$$\mathbf{E}[I_t^{H_b,b}] = \begin{cases} y_t^{D,b} & \forall 1 \leq t \leq L+2; \\ y_t^{D,b} + \sum_{s=2}^{t-1-L} k_{s-1}^b & \forall L+3 \leq t \leq T. \end{cases}$$

The proofs of Lemmas II.6 and II.7 are by induction and provided in the Appendix (the identities in Lemma II.6 are used to prove Lemma II.7).

Step 2

We derive an upper bound for $C^{H_b,b} - D^{*,b}$. By definition,

$$C^{H_b,b} = \mathbf{E} \left[\sum_{t=1}^T c \cdot x_{t-L}^{H_b,b} + \sum_{t=1}^T h \cdot (I_t^{H_b,b} + x_{t-L}^{H_b,b} - D_t)^+ + \sum_{t=1}^T p \cdot (D_t - x_{t-L}^{H_b,b} - I_t^{H_b,b})^+ \right] \quad (2.7)$$

For the first summation in (2.7), by definition of k_t^b , we immediately have:

$$\mathbf{E} \left[\sum_{t=1}^T c \cdot x_{t-L}^{H_b,b} \right] = \sum_{t=1}^T c \cdot x_{t-L}^{D,b} + \sum_{t=2}^{T-L} c \cdot k_{t-1}^b. \quad (2.8)$$

As for the second summation in (2.7), note that

$$\begin{aligned} & \mathbf{E} \left[\sum_{t=1}^T (I_t^{H_b,b} + x_{t-L}^{H_b,b} - D_t)^+ \right] \\ &= \mathbf{E} \left[\sum_{t=1}^T (I_t^{H_b,b} + x_{t-L}^{H_b,b} - D_t) + \sum_{t=1}^T (D_t - x_{t-L}^{H_b,b} - I_t^{H_b,b})^+ \right] \\ &= \sum_{t=1}^T (y_t^{D,b} + x_{t-L}^{D,b} - \mu_t) + \sum_{t=L+3}^T \sum_{s=2}^{t-1-L} k_{s-1}^b + \sum_{t=2}^{T-L} k_{t-1}^b + \mathbf{E} \left[\sum_{t=1}^T (D_t - x_{t-L}^{H_b,b} - I_t^{H_b,b})^+ \right] \\ &= \sum_{t=1}^T (y_t^{D,b} + x_{t-L}^{D,b} - \mu_t) + \sum_{t=2}^{T-L} (T-L-t+1) \cdot k_{t-1}^b + \mathbf{E} \left[\sum_{t=1}^T (D_t - x_{t-L}^{H_b,b} - I_t^{H_b,b})^+ \right], \end{aligned} \quad (2.9)$$

where the second equality follows from the definition of k_t^b and Lemma II.7. By (2.3) and (2.4), the terms $\mathbf{E}[(D_t - x_{t-L}^{H_b,b} - I_t^{H_b,b})^+]$ can be bounded as follows:

$$\mathbf{E}[(D_t - x_{t-L}^{H_b,b} - I_t^{H_b,b})^+] \leq \mathbf{E} \left[\left(\sum_{s=1}^t D_s - \sum_{s=1}^t \mu_s - \beta_1^t(\alpha) \right)^+ \right] \quad \forall 1 \leq t \leq L+1, \quad (2.10)$$

$$\mathbf{E}[(D_t - x_{t-L}^{H_b,b} - I_t^{H_b,b})^+] \leq \mathbf{E} \left[\left(\sum_{s=t-L}^t D_s - \sum_{s=t-L}^t \mu_s - \beta_{t-L}^{L+1}(\alpha) \right)^+ \right] \quad \forall L+2 \leq t \leq T. \quad (2.11)$$

Combining (2.8), (2.9), (2.10), and (2.11) and noting that, by definition,

$$D^{*,b} \geq \sum_{t=1}^T c \cdot x_{t-L}^{D,b} + \sum_{t=1}^T h \cdot y_{t+1}^{D,b} = \sum_{t=1}^T c \cdot x_{t-L}^{D,b} + \sum_{t=1}^T h \cdot (y_t^{D,b} + x_{t-L}^{D,b} - \mu_t)$$

immediately yields

$$\begin{aligned} C^{H_b,b} - D^{*,b} &\leq \sum_{t=2}^{T-L} [c + h \cdot (T - L - t + 1)] \cdot \mathbf{E} \left[x_t^{H_b,b}(\alpha) - x_t^{D,b}(\alpha) \right] \\ &\quad + (h + p) \cdot \sum_{t=1}^{L+1} \mathbf{E} \left[\left(\sum_{s=1}^t D_s - \sum_{s=1}^t \mu_s - \beta_1^t(\alpha) \right)^+ \right] \\ &\quad + (h + p) \cdot \sum_{t=L+2}^T \mathbf{E} \left[\left(\sum_{s=t-L}^t D_s - \sum_{s=t-L}^t \mu_s - \beta_{t-L}^{L+1}(\alpha) \right)^+ \right]. \end{aligned}$$

Step 3

To derive a lower bound of $C^{*,b}$, by Lemma II.2, $C^{*,b} \geq D^{*,b}$. By the constraints in (2.5), we know that $\sum_{t=1}^T x_{t-L} \geq \sum_{t=1}^T \mu_t + \beta_{T-L}^{L+1}(\alpha)$, $y_{t+1} \geq \beta_1^t(\alpha)$ for $1 \leq t \leq L+1$, $y_{t+1} \geq \beta_{t-L}^{L+1}(\alpha)$ for $L+2 \leq t \leq T$, and $p_t \geq 0$ for all t . So, we can bound:

$$C^{*,b} \geq D^{*,b} \geq c \cdot \left[\sum_{t=1}^T \mu_t + \beta_{T-L}^{L+1}(\alpha) \right] + h \cdot \left[\sum_{t=1}^{L+1} \beta_1^t(\alpha) + \sum_{t=L+2}^T \beta_{t-L}^{L+1}(\alpha) \right].$$

This completes the proof of Theorem II.5. ■

2.3.3 Numerical Results

In this section, we numerically investigate the performance of H_b . We assume that demands are non-i.i.d. (specifically, demands are cyclic in the sense that their means follow the following pattern: 5, 6, 7, 8, 9, 10, 9, 8, 7, 6, 5, 6, ...) and report the percentage performance gaps of H_b with respect to the LP lower bound. In Tables 2.1 and 2.2, demand follows Poisson or Normal Distribution respectively (similar to Bitran and Yanasse [19]) and is truncated between $[0, \bar{D}]$. In Table 2.1, the mean (and standard deviation) of the untruncated demand varies from 5 to 10 and \bar{D} is set to be 23. In Table 2.2, the mean in the untruncated demand varies from 5 to 10, the standard deviation is set to be 2, and \bar{D} is set to be 16.

In the numerical experiment, we set value of the other problem parameters either similar to the practical settings or similar to that used in the literature (Zipkin [95]). Specifically, the length of the time horizon varies from 10 to 40 (the practical planning horizon is usually no longer than several weeks); the length of lead time varies from 1 to 4; the ordering cost is set to 0; the holding cost is normalized to 1; and the penalty cost varies from 2 to 20. In addition, we also vary service level requirement from 0.9 to 0.999. As noted in the introduction, we are primarily interested in the setting with a high service level. Thus, we set the smallest service level to be 0.9.

In both cases (Normal and Poisson demands), the performance gap of H_b converges to 0 when the service-level increases. Interestingly, we observe that the percentage gap is not very sensitive with T , especially when the penalty cost is high. Thus, although our theoretical bound in Theorem II.5 suggests that the performance gap might significantly worsen when T is large, the actual performance of H_b is very promising. As for the impact of the penalty cost on performance, the gaps grow almost linearly with the

		T=10				T=20				T=40			
α		1	2	3	4	1	2	3	4	1	2	3	4
p=2	0.001	0.04	0.04	0.03	0.03	0.05	0.09	0.12	0.15	0.06	0.10	0.12	0.17
	0.01	0.49	0.52	0.47	0.47	0.52	0.54	0.51	0.60	0.54	0.56	0.56	0.62
	0.05	4.12	4.26	4.23	4.23	4.26	4.28	4.33	4.41	4.21	4.28	4.40	4.47
	0.1	12.49	12.30	11.06	10.95	12.37	12.20	11.64	11.66	12.77	12.18	11.78	11.66
p=4	0.001	0.06	0.07	0.05	0.05	0.06	0.11	0.12	0.12	0.09	0.12	0.16	0.17
	0.01	0.80	0.84	0.80	0.78	0.86	0.88	0.87	0.91	0.89	0.87	0.85	0.93
	0.05	6.88	7.04	7.17	6.99	7.15	7.05	7.20	7.34	7.13	7.07	7.30	7.44
	0.1	20.71	20.17	18.71	18.57	20.67	20.19	19.20	19.45	21.08	20.34	19.49	19.45
p=6	0.001	0.08	0.10	0.07	0.08	0.09	0.11	0.14	0.13	0.10	0.15	0.18	0.18
	0.01	1.14	1.14	1.17	1.20	1.19	1.17	1.16	1.15	1.22	1.19	1.22	1.31
	0.05	9.66	9.92	9.81	9.78	9.61	9.87	10.22	9.86	10.00	9.95	10.03	10.15
	0.1	29.01	28.48	26.43	25.48	28.76	28.51	27.04	26.95	29.41	28.26	27.30	27.00
p=8	0.001	0.10	0.11	0.11	0.09	0.14	0.15	0.18	0.19	0.13	0.18	0.18	0.21
	0.01	1.45	1.46	1.51	1.46	1.51	1.50	1.44	1.58	1.56	1.51	1.47	1.55
	0.05	12.50	12.57	12.77	12.69	12.44	12.73	12.91	12.84	12.79	12.71	12.88	13.09
	0.1	37.45	36.81	33.80	33.14	37.20	36.35	34.59	34.54	38.17	36.21	35.06	34.91

Table 2.1: Percentage Performance Gap (in %) with Poisson Distributed Demand

		T=10				T=20				T=40			
		lead time				lead time				lead time			
α		1	2	3	4	1	2	3	4	1	2	3	4
p=2	0.001	0.02	0.02	0.03	0.03	0.03	0.03	0.03	0.03	0.03	0.03	0.03	0.03
	0.01	0.43	0.44	0.42	0.43	0.43	0.44	0.43	0.44	0.43	0.44	0.44	0.44
	0.05	3.77	3.77	3.81	3.79	3.76	3.81	3.82	3.83	3.81	3.80	3.81	3.84
	0.1	11.01	10.98	11.12	11.04	11.03	11.03	11.05	11.13	11.08	11.10	11.10	11.10
p=4	0.001	0.04	0.05	0.04	0.04	0.05	0.05	0.05	0.05	0.05	0.05	0.05	0.05
	0.01	0.72	0.73	0.72	0.73	0.72	0.72	0.73	0.73	0.73	0.73	0.73	0.73
	0.05	6.33	6.38	6.35	6.34	6.30	6.36	6.37	6.35	6.32	6.34	6.38	6.35
	0.1	18.38	18.41	18.27	18.38	18.37	18.40	18.49	18.48	18.46	18.46	18.42	18.50
p=6	0.001	0.06	0.06	0.06	0.06	0.06	0.07	0.07	0.07	0.07	0.07	0.07	0.07
	0.01	0.99	1.01	1.00	1.00	1.01	1.01	1.02	1.04	1.01	1.02	1.02	1.03
	0.05	8.87	8.94	8.87	8.85	8.82	8.89	8.90	8.84	8.85	8.90	8.89	8.86
	0.1	25.68	25.79	25.78	25.98	25.78	25.93	25.93	25.92	25.76	25.80	25.87	25.89
p=8	0.001	0.07	0.08	0.08	0.08	0.08	0.09	0.08	0.08	0.08	0.09	0.09	0.09
	0.01	1.26	1.29	1.28	1.29	1.27	1.29	1.31	1.31	1.30	1.31	1.32	1.30
	0.05	11.41	11.37	11.37	11.32	11.36	11.31	11.33	11.36	11.35	11.43	11.42	11.43
	0.1	33.00	33.19	32.98	33.06	33.12	33.27	33.35	33.22	33.12	33.26	33.19	33.26

Table 2.2: Percentage Performance Gap (in %) with Normal Distributed Demand

penalty cost, which is not surprising and is predicted by the bound in Theorem II.5. The intuition is as follows. As the penalty cost increases, the optimal order quantity should also increase. However, since our deterministic approximation is mainly designed to handle the service-level constraints, as can be seen from Lemma 1, its solution is, unfortunately, not affected by the penalty cost. This suggests a limitation in applying deterministic approximation to settings with large penalty costs. We leave it as an interesting open question how to approximate an inventory problem with both service level constraints and a large penalty cost.

2.4 Lost-Sales Inventory System

We now focus on the lost-sales system. We first describe our heuristic control and discuss its theoretical performance in Section 4.1, and then we provide an outline of the proof of our main result (Theorem II.9) in Section 4.2.

2.4.1 Proposed Heuristic Control and Its Performance

We call our heuristic control for the lost-sales system simply as H_ℓ . Given the service level $1 - \alpha$, in period t , H_ℓ places a new order $x_t^{H_\ell, \ell}(\alpha)$ as follows:

$$x_t^{H_\ell, \ell}(\alpha) = \begin{cases} x_t^{D,b}(\alpha) & \forall 1 - L \leq t \leq 1; \\ \left(y_t^{D,b}(\alpha) + \sum_{s=t-L}^t x_s^{D,b}(\alpha) - I_t^{H_\ell, \ell} - \sum_{s=t-L}^{t-1} x_s^{H_\ell, \ell}(\alpha) \right)^+ & \forall 2 \leq t \leq T - L. \end{cases}$$

Note that H_ℓ uses the same order-up-to level as control H_b defined in (2.6) in Section 2.3.1. Moreover, by construction, H_ℓ is a feasible control to the stochastic lost-sale system. We state this formally below.

Lemma II.8. $H_\ell \in \Pi^\ell$.

Let $\phi(T, \alpha) = \sum_{t=1}^{T-1} \sum_{s=1}^{t-1} \theta_s(\alpha)$ where $\theta_t(\alpha)$ is defined as follows:

$$\theta_t(\alpha) = \begin{cases} t \cdot \bar{D} - \sum_{s=1}^t \mu_s - \beta_1^t(\alpha) & \forall 1 \leq t \leq L+1; \\ (L+1) \cdot \bar{D} - \sum_{s=t-L}^t \mu_s - \beta_{t-L}^{L+1}(\alpha) & \forall L+2 \leq t \leq T. \end{cases}$$

Note that, by definition, $\theta_t(\alpha) \rightarrow 0$ as $\alpha \rightarrow 0$ for all t .

We present our main result for the lost-sale system below.

Theorem II.9. *We can bound:*

$$\begin{aligned} C^{H_\ell, \ell}(\alpha) - C^{*, \ell}(\alpha) &\leq (c+h) \cdot \phi(T, \alpha) + (c+L \cdot p) \cdot T \cdot \alpha \cdot \bar{D} \\ &\quad + (c+h) \cdot \sum_{t=2}^{T-L} (T-L-t+1) \cdot \mathbf{E} \left[x_t^{H_b, b}(\alpha) - x_t^{D, b}(\alpha) \right] \\ &\quad + (h+p) \cdot \sum_{t=1}^{L+1} (T-t+1) \cdot \mathbf{E} \left[\left(\sum_{s=1}^t D_s - \sum_{s=1}^t \mu_s - \beta_1^t(\alpha) \right)^+ \right] \\ &\quad + (h+p) \cdot \sum_{t=L+2}^T (T-t+1) \cdot \mathbf{E} \left[\left(\sum_{s=t-L}^t D_s - \sum_{s=t-L}^t \mu_s - \beta_{t-L}^{L+1}(\alpha) \right)^+ \right]. \end{aligned}$$

Moreover, we also have:

$$\begin{aligned} C^{*, \ell}(\alpha) &\geq c \cdot \left[\sum_{t=1}^T \mu_t + \beta_{T-L}^{L+1}(\alpha) \right] + h \cdot \left[\sum_{t=1}^{L+1} \beta_1^t(\alpha) + \sum_{t=L+2}^T \beta_{t-L}^{L+1}(\alpha) \right] \\ &\quad - 2 \cdot (c+h) \cdot \phi(T, \alpha) - (c+L \cdot p) \cdot T \cdot \alpha \cdot \bar{D}. \end{aligned}$$

Similar to the discussions following Theorem II.5, it is not difficult to see that the bound in Theorem II.9 contains $O(T^2)$ terms, each of which converges to 0 as $\alpha \rightarrow 0$, even for the case when demands are i.i.d. It remains unclear to us whether

a deterministic approximation can be used to construct a policy whose bound can be written as a sum of $O(T)$ terms each of whom converges to 0 as $\alpha \rightarrow 0$, at least for the case of i.i.d demands. We leave it as an interesting open question.

2.4.2 Proof of Theorem II.9

As in the proof of Theorem II.5, with an exception of a few notations, we will often suppress notational dependency on α . Note that $C^{H_\ell, \ell} - C^{*, \ell}$ can be written as a sum of two terms:

$$C^{H_\ell, \ell} - C^{*, \ell} = [C^{H_\ell, \ell} - D^{*, b}] + [D^{*, b} - C^{*, \ell}].$$

In what follows, we will divide the proof of Theorem II.9 into two major parts: in part 1 (Section 2.4.2.1), we derive an upper bound for $C^{H_\ell, \ell} - D^{*, b}$; in part 2 (Section 2.4.2.2), we derive an upper bound for $D^{*, b} - C^{*, \ell}$. A lower bound for $C^{*, \ell}$ is also discussed in Section 2.4.2.2.

2.4.2.1 An Upper Bound for $C^{H_\ell, \ell} - D^{*, b}$.

We state our main proposition for part 1 below.

Proposition II.10. We can bound:

$$\begin{aligned} C^{H_\ell, \ell} - D^{*, b} &\leq (c + h) \cdot \sum_{t=2}^{T-L} (T - L - t + 1) \cdot k_{t-1}^b \\ &\quad + (h + p) \cdot \sum_{t=1}^{L+1} (T - t + 1) \cdot \mathbf{E} \left[\left(\sum_{s=1}^t D_s - \sum_{s=1}^t \mu_s - \beta_1^t(\alpha) \right)^+ \right] \\ &\quad + (h + p) \cdot \sum_{t=L+2}^T (T - t + 1) \cdot \mathbf{E} \left[\left(\sum_{s=t-L}^t D_s - \sum_{s=t-L}^t \mu_s - \beta_{t-L}^{L+1}(\alpha) \right)^+ \right]. \end{aligned}$$

Below we state a lemma that will be useful for proving Proposition II.10.

Lemma II.11. $x_t^{H_\ell, \ell} \leq x_t^{H_b, b} \quad \forall t \geq 1 - L.$

The proof of Lemma II.11 can be found in the Appendix. We now proceed to prove Proposition II.10 in two steps. In Step 1, we provide upper bounds for $\mathbf{E}[(D_t - x_{t-L}^{H_\ell, \ell} - I_t^{H_\ell, \ell})^+]$; in Step 2, we use the bound derived in Step 1, together with the identities derived in Lemmas II.11 to bound $C^{H_\ell, \ell}$; the result in Proposition II.10 immediately follows by subtracting $D^{*,b}$ from the bound derived in Step 2.

Step 1

We claim that

$$\begin{aligned} \mathbf{E}[(D_t - x_{t-L}^{H_\ell, \ell} - I_t^{H_\ell, \ell})^+] &\leq \mathbf{E} \left[\left(\sum_{s=1}^t D_s - \sum_{s=1}^t \mu_s - \beta_1^t(\alpha) \right)^+ \right] && \forall 1 \leq t \leq L+1, \\ \mathbf{E}[(D_t - x_{t-L}^{H_\ell, \ell} - I_t^{H_\ell, \ell})^+] &\leq \mathbf{E} \left[\left(\sum_{s=t-L}^t D_s - \sum_{s=t-L}^t \mu_s - \beta_{t-L}^{L+1}(\alpha) \right)^+ \right] && \forall L+2 \leq t \leq T. \end{aligned}$$

As $x_s^{H_\ell, \ell} = x_s^{D,b}$ for $s \leq 1$, for $1 \leq t \leq L+1$, we immediately have

$$\begin{aligned} \mathbf{E}[(D_t - x_{t-L}^{H_\ell, \ell} - I_t^{H_\ell, \ell})^+] &\leq \mathbf{E} \left[\left(\sum_{s=1}^t D_s - \sum_{s=1}^t x_{s-L}^{D,b} \right)^+ \right] \\ &\leq \mathbf{E} \left[\left(\sum_{s=1}^t D_s - \sum_{s=1}^t \mu_s - \beta_1^t(\alpha) \right)^+ \right]. \end{aligned}$$

The first inequality follows from the fact that, given the same order quantities, the inventory level in the lost-sale system cannot be smaller than the inventory level in the backorder system (i.e., $I_t^{H_\ell, \ell} \geq I_t^{H_b, b}$ for $1 \leq t \leq L+1$). The second inequality follows from the constraints in (2.5), i.e.,

$$\sum_{s=1}^t x_{s-L}^{D,b} - \sum_{s=1}^t \mu_s \geq \beta_1^t(\alpha).$$

As for periods $t \geq L + 2$, note that

$$\begin{aligned}
I_t^{H_\ell, \ell} + x_{t-L}^{H_\ell, \ell} - D_t &\geq I_{t-L}^{H_\ell, \ell} + \sum_{s=t-L}^t x_{s-L}^{H_\ell, \ell} - \sum_{s=t-L}^t D_s \\
&\geq y_{t-L}^{D,b} + \sum_{s=t-L}^t x_{s-L}^{D,b} - \sum_{s=t-L}^t D_s,
\end{aligned} \tag{2.12}$$

where the second inequality follows because $x_k^{H_\ell, \ell} \geq y_k^{D,b} + \sum_{s=k-L}^k x_s^{D,b} - I_k^{H_\ell, \ell} - \sum_{s=k-L}^{k-1} x_s^{H_\ell, \ell}$ for $k \geq 2$ (in particular, $x_{t-L}^{H_\ell, \ell} \geq y_{t-L}^{D,b} + \sum_{s=t-2L}^{t-L} x_s^{D,b} - I_{t-L}^{H_\ell, \ell} - \sum_{s=t-2L}^{t-L-1} x_s^{H_\ell, \ell}$, which implies (2.12)). By the constraints in (2.5) again, for $L+2 \leq t \leq T$, we therefore have:

$$\begin{aligned}
\mathbf{E}[(D_t - x_{t-L}^{H_\ell, \ell} - I_t^{H_\ell, \ell})^+] &\leq \mathbf{E} \left[\left(\sum_{s=t-L}^t D_s - \sum_{s=t-L}^t x_{s-L}^{D,b} - y_{t-L}^{D,b} \right)^+ \right] \\
&\leq \mathbf{E} \left[\left(\sum_{s=t-L}^t D_s - \sum_{s=t-L}^t \mu_s - \beta_{t-L}^{L+1}(\alpha) \right)^+ \right].
\end{aligned}$$

Step 2

We now explicitly bound the expected total costs incurred by H_ℓ in the lost-sales system:

$$C^{H_\ell, \ell} = \mathbf{E} \left[\sum_{t=1}^T c \cdot x_{t-L}^{H_\ell, \ell} + \sum_{t=1}^T h \cdot (I_t^{H_\ell, \ell} + x_{t-L}^{H_\ell, \ell} - D_t)^+ + \sum_{t=1}^T p \cdot (D_t - I_t^{H_\ell, \ell} - x_{t-L}^{H_\ell, \ell})^+ \right] \tag{2.13}$$

We first do the following transformation on the total order quantities. Note that for any $\pi \in \Pi^\ell$, we have the following identity (2.14) due to the fact that the inventory on-hand in the end of period t (the left side of (2.14)) equals the total supply over periods $[1, t]$ (the first term on the right side of (2.14)) minus the total demand over periods $[1, t]$ (the second term on the right side of (2.14)) plus the total lost-sales over

periods $[1, t]$ (the third term on the right side of (2.14)):

$$(I_t^{\pi, \ell} + x_{t-L}^{\pi, \ell} - D_t)^+ = \sum_{s=1}^t x_{s-L}^{\pi, \ell} - \sum_{s=1}^t D_s + \sum_{s=1}^t (D_s - I_s^{\pi, \ell} - x_{s-L}^{\pi, \ell})^+. \quad (2.14)$$

It follows that

$$\sum_{s=1}^T x_{s-L}^{H_\ell, \ell} = (I_T^{H_\ell, \ell} + x_{T-L}^{H_\ell, \ell} - D_T)^+ + \sum_{s=1}^T D_s - \sum_{s=1}^T (D_s - I_s^{H_\ell, \ell} - x_{s-L}^{H_\ell, \ell})^+.$$

Using this, (2.13) can be re-written as

$$\begin{aligned} C^{H_\ell, \ell} &= \mathbf{E} \left[h \cdot \sum_{t=1}^{T-1} (I_t^{H_\ell, \ell} + x_{t-L}^{H_\ell, \ell} - D_t)^+ + (c+h) \cdot (I_T^{H_\ell, \ell} + x_{T-L}^{H_\ell, \ell} - D_T)^+ \right. \\ &\quad \left. + c \cdot \sum_{t=1}^T D_t + (p-c) \cdot \sum_{t=1}^T (D_t - I_t^{H_\ell, \ell} - x_{t-L}^{H_\ell, \ell})^+ \right]. \end{aligned} \quad (2.15)$$

The first summation in (2.15) can be bounded as follows:

$$\begin{aligned} &\mathbf{E} \left[\sum_{t=1}^{T-1} (I_t^{H_\ell, \ell} + x_{t-L}^{H_\ell, \ell} - D_t)^+ \right] \\ &= \mathbf{E} \left[\sum_{t=1}^{T-1} \left(\sum_{s=1}^t x_{s-L}^{H_\ell, \ell} - \sum_{s=1}^t D_s + \sum_{s=1}^t (D_s - x_{s-L}^{H_\ell, \ell} - I_s^{H_\ell, \ell})^+ \right) \right] \\ &\leq \mathbf{E} \left[\sum_{t=1}^{T-1} \left(\sum_{s=1}^t x_{s-L}^{H_b, b} - \sum_{s=1}^t D_s + \sum_{s=1}^t (D_s - x_{s-L}^{H_\ell, \ell} - I_s^{H_\ell, \ell})^+ \right) \right] \\ &= \sum_{t=1}^{T-1} \left\{ \sum_{s=1}^t x_{s-L}^{D, b} + \sum_{s=2}^{t-L} k_{s-1}^b - \sum_{s=1}^t \mu_s + \mathbf{E} \left[\sum_{s=1}^t (D_s - x_{s-L}^{H_\ell, \ell} - I_s^{H_\ell, \ell})^+ \right] \right\} \\ &= \sum_{t=1}^{T-1} (y_t^{D, b} + x_{t-L}^{D, b} - \mu_t) + \sum_{t=2}^{T-L-1} (T-L-t) \cdot k_{t-1}^b + \mathbf{E} \left[\sum_{t=1}^{T-1} (T-t) \cdot (D_t - x_{t-L}^{H_\ell, \ell} - I_t^{H_\ell, \ell})^+ \right], \end{aligned}$$

where the first equality again follows by repeatedly applying identity $z^+ = z + (-z)^+$

and the first inequality follows from Lemma II.11. Similarly,

$$\mathbf{E} \left[(I_T^{H_\ell, \ell} + x_{T-L}^{H_\ell, \ell} - D_T)^+ \right] \leq (y_T^{D,b} + x_{T-L}^{D,b} - \mu_T) + \sum_{t=2}^{T-L} k_{t-1}^b + \mathbf{E} \left[\sum_{t=1}^T (D_t - x_{t-L}^{H_\ell, \ell} - I_t^{H_\ell, \ell})^+ \right].$$

Now, by definition,

$$\begin{aligned} D^{*,b} &\geq \sum_{t=1}^T c \cdot x_{t-L}^{D,b} + \sum_{t=1}^T h \cdot (y_t^{D,b} + x_{t-L}^{D,b} - \mu_t) \\ &= \sum_{t=1}^{T-1} h \cdot (y_t^{D,b} + x_{t-L}^{D,b} - \mu_t) + (c+h) \cdot (y_T^{D,b} + x_{T-L}^{D,b} - \mu_T) + c \cdot \sum_{s=1}^T \mu_s. \end{aligned}$$

Thus, we can bound (2.15) as follows:

$$\begin{aligned} C^{H_\ell, \ell} - D^{*,b} &\leq h \cdot \mathbf{E} \left[\sum_{t=1}^{T-1} (T-t) \cdot (D_t - x_{t-L}^{H_\ell, \ell} - I_t^{H_\ell, \ell})^+ \right] \\ &\quad + (c+h) \cdot \mathbf{E} \left[\sum_{t=1}^T (D_t - x_{t-L}^{H_\ell, \ell} - I_t^{H_\ell, \ell})^+ \right] \\ &\quad + (p-c) \cdot \mathbf{E} \left[\sum_{t=1}^T (D_t - I_t^{H_\ell, \ell} - x_{t-L}^{H_\ell, \ell})^+ \right] \\ &\quad + (c+h) \cdot \sum_{t=2}^{T-L} (T-L-t+1) \cdot k_{t-1}^b \\ &\leq (h+p) \cdot \mathbf{E} \left[\sum_{t=1}^T (T-t+1) \cdot (D_t - x_{t-L}^{H_\ell, \ell} - I_t^{H_\ell, \ell})^+ \right] \\ &\quad + (c+h) \cdot \sum_{t=2}^{T-L} (T-L-t+1) \cdot k_{t-1}^b. \end{aligned}$$

This completes the proof.

2.4.2.2 An Upper Bound for $D^{*,b} - C^{*,\ell}$ and a Lower Bound for $C^{*,\ell}$.

We state a proposition.

Proposition II.12. We can bound:

$$D^{*,b} - C^{*,\ell} \leq (c + h) \cdot \phi(T, \alpha) + (c + L \cdot p) \cdot T \cdot \alpha \cdot \bar{D}.$$

Moreover, we also have:

$$\begin{aligned} C^{*,\ell} \geq & c \cdot \left[\sum_{t=1}^T \mu_t + \beta_{T-L}^{L+1}(\alpha) \right] + h \cdot \left[\sum_{t=1}^{L+1} \beta_1^t(\alpha) + \sum_{t=L+2}^T \beta_{t-L}^{L+1}(\alpha) \right] \\ & - 2 \cdot (c + h) \cdot \phi(T, \alpha) - (c + L \cdot p) \cdot T \cdot \alpha \cdot \bar{D}. \end{aligned}$$

Unlike the quantities $C^{H_\ell, \ell}$ and $D^{*,b}$ in subsection 2.4.2.1 that are not too difficult to compare by exploiting the relationships among the key variables in the lost-sales, backorder, and deterministic systems, which can be explicitly derived (cf. Proposition II.10), the quantities $D^{*,b}$ and $C^{*,\ell}$ are not easily comparable since the optimal control for the lost-sales system is not known. One natural approach would be to first decompose $D^{*,b} - C^{*,\ell}$ into a sum of two terms, i.e., $[D^{*,b} - C^{*,b}] + [C^{*,b} - C^{*,\ell}]$, and then bound each of the terms separately. However, note that while $D^{*,b} - C^{*,b} \leq 0$ (by Lemma II.2, so this term can be ignored), the term $C^{*,b} - C^{*,\ell}$ is not easy to bound directly. To bypass this difficulty, in proving the upper bound in Proposition II.12, we will construct an alternative backorder system \tilde{b} whose optimal expected total costs is only slightly larger than the optimal expected total costs under the lost-sales system. Note that although not the same, this comparison shares the same spirit as the comparison between the backorder and lost-sales systems in the canonical cost-based model; e.g., Janakiraman, Seshadri and Shanthikumar [48].

Define $C^{*,\tilde{b}}$ as follows:

$$C^{*,\tilde{b}} = \min_{\pi \in \Pi^{\tilde{b}}} \mathbf{E} \left[\sum_{t=1}^T c \cdot x_{t-L}^{\pi, \tilde{b}} + \sum_{t=1}^T h \cdot (I_t^{\pi, \tilde{b}} + x_{t-L}^{\pi, \tilde{b}} - D_t)^+ + \sum_{t=1}^T p \cdot (D_t - I_t^{\pi, \tilde{b}} - x_{t-L}^{\pi, \tilde{b}})^+ \right],$$

where the set of feasible controls $\Pi^{\tilde{b}}$ is defined as:

$$\Pi^{\tilde{b}} = \left\{ \pi : x_t^{\pi, \tilde{b}} \geq 0, P \left(I_t^{\pi, \tilde{b}} + x_{t-L}^{\pi, \tilde{b}} - D_t \geq - \sum_{s=1}^{t-1} \theta_s(\alpha) \right) \geq 1 - \alpha, \forall 1 \leq t \leq L + 1 \text{ and} \right. \\ \left. P \left(I_t^{\pi, \tilde{b}} + x_{t-L}^{\pi, \tilde{b}} - D_t \geq - \sum_{s=t-L}^{t-1} \theta_s(\alpha) \middle| \mathfrak{S}_{t-L-1}^{\pi, \tilde{b}} \right) \geq 1 - \alpha, \forall L + 2 \leq t \leq T \right\},$$

and $I_{t+1}^{\pi, \tilde{b}} = I_t^{\pi, \tilde{b}} + x_{t-L}^{\pi, \tilde{b}} - D_t$ for all t .

Define $\gamma_t(\alpha)$ and $w_t(\alpha)$ as follows:

- For $1 \leq t \leq L + 1$, $\gamma_t(\alpha)$ is the smallest γ that satisfies

$$P \left(\gamma + \sum_{s=1}^t \mu_s - \sum_{s=1}^t D_s \geq - \sum_{s=1}^{t-1} \theta_s(\alpha) \right) \geq 1 - \alpha.$$

- For $L + 2 \leq t \leq T$, $w_t(\alpha)$ is the smallest w that satisfies

$$P \left(w + \sum_{s=t-L}^t \mu_s - \sum_{s=t-L}^t D_s \geq - \sum_{s=t-L}^{t-1} \theta_s(\alpha) \right) \geq 1 - \alpha.$$

The probabilistic service level constraints in $\Pi^{\tilde{b}}$ can be equivalently written as follows:

$$I_1^{\pi, \tilde{b}} + \sum_{s=1}^t x_{s-L}^{\pi, \tilde{b}} - \sum_{s=1}^t \mu_s \geq \gamma_t(\alpha) \quad \forall 1 \leq t \leq L + 1, \\ I_{t-L}^{\pi, \tilde{b}} + \sum_{s=t-L}^t x_{s-L}^{\pi, \tilde{b}} - \sum_{s=t-L}^t \mu_s \geq w_t(\alpha) \quad \forall L + 2 \leq t \leq T.$$

The following deterministic optimization is the analogue of (2.5) for backorder system

\tilde{b} :

$$\begin{aligned}
D^{*,\tilde{b}} = & \min_{x,y} \sum_{t=1}^T c \cdot x_{t-L} + \sum_{t=1}^T h \cdot z_{t+1} + \sum_{t=1}^T p \cdot m_{t+1} & (2.16) \\
\text{s.t. } & y_1 = 0 \\
& y_t = y_{t-1} + x_{t-1-L} - \mu_{t-1} \quad \forall 2 \leq t \leq T+1 \\
& z_t \geq y_t \quad \forall 2 \leq t \leq T+1 \\
& m_t \geq -y_t \quad \forall 2 \leq t \leq T+1 \\
& \sum_{s=1}^t x_{s-L} - \sum_{s=1}^t \mu_s \geq \gamma_t(\alpha) \quad \forall 1 \leq t \leq L+1 \\
& \sum_{s=1}^t x_{s-L} - \sum_{s=1}^t \mu_s \geq w_t(\alpha) \quad \forall L+2 \leq t \leq T \\
& x_t, z_t, m_t \geq 0 \quad \forall 1-L \leq t \leq T-L.
\end{aligned}$$

The following relations between $\beta_k^t(\alpha)$, $\gamma_t(\alpha)$, and $w_t(\alpha)$ are useful. By the definitions of $\beta_t^k(\alpha)$, $\gamma_t(\alpha)$, $w_t(\alpha)$, we can write:

$$\begin{aligned}
\gamma_t(\alpha) &= \beta_1^t(\alpha) - \sum_{s=1}^{t-1} \theta_s(\alpha) \quad \forall 1 \leq t \leq L+1 \\
w_t(\alpha) &= \beta_{t-L}^{L+1}(\alpha) - \sum_{s=1}^{t-1} \theta_s(\alpha) \quad \forall L+2 \leq t \leq T
\end{aligned}$$

Since $\theta_s(\alpha) \geq 0$, we have:

$$\begin{aligned}
\gamma_t(\alpha) &\leq \beta_1^t(\alpha) \quad \forall 1 \leq t \leq L+1, \\
w_t(\alpha) &\leq \beta_{t-L}^{L+1}(\alpha) \quad \forall L+2 \leq t \leq T.
\end{aligned}$$

In addition, noting that $\gamma_t(0) = \beta_1^t(0)$ and $w_t(0) = \beta_{t-L}^{L+1}(0)$, we have

$$\gamma_t(\alpha) \rightarrow \beta_1^t(\alpha), \quad w_t(\alpha) \rightarrow \beta_{t-L}^{L+1}(\alpha) \quad \text{as } \alpha \rightarrow 0.$$

Intuitively, one can interpret $\gamma_t(\alpha)$ and $\omega_t(\alpha)$ as the adjusted “ $(1 - \alpha)$ percentiles”, which are adjusted by the maximum possible accumulated “lost” reflected in θ 's.

Below, we state several lemmas that will be useful for proving Proposition II.12; their proofs can be found in the Appendix.

Lemma II.13. $D^{*,\tilde{b}} \leq C^{*,\tilde{b}}$.

Lemma II.14. *Under any feasible control $\pi \in \Pi^\ell$, the following holds for all $t \geq 1$:*

$$I_t^{\pi,\ell} + x_{t-L}^{\pi,\ell} - D_t = - \sum_{s=1}^t \Upsilon_s + \max \left\{ 0, \Upsilon_1, \sum_{s=1}^2 \Upsilon_s, \dots, \sum_{s=1}^{t-1} \Upsilon_s \right\}$$

where $\Upsilon_s = D_s - x_{s-L}^{\pi,\ell}$.

Lemma II.15. *For any $\pi \in \Pi^\ell$, if we apply the same ordering decision x_t^π at period t in the backorder system \tilde{b} as if the inventory levels evolve according to a lost-sales system, the resulting sequence of ordering decisions satisfies the probabilistic service level constraints in $\Pi^{\tilde{b}}$, i.e., $\Pi^\ell \subseteq \Pi^{\tilde{b}}$.*

Note that Lemma II.15 does not imply $C^{*,\tilde{b}} - C^{*,\ell} \leq 0$, as the backorder systems \tilde{b} and the lost-sales system ℓ have different inventory evolution dynamics.

Lemma II.16. $C^{*,\tilde{b}} - C^{*,\ell} \leq (c + L \cdot p) \cdot T \cdot \alpha \cdot \bar{D}$.

Lemma II.17. *We can bound $D^{*,b} - D^{*,\tilde{b}}$ as follows:*

$$D^{*,b} - D^{*,\tilde{b}} \leq (c + h) \cdot \phi(T, \alpha).$$

To prove the upper bound in Proposition II.12, note that we can decompose $D^{*,b} - C^{*,\ell}$ as follows:

$$D^{*,b} - C^{*,\ell} = [D^{*,b} - D^{*,\tilde{b}}] + [D^{*,\tilde{b}} - C^{*,\tilde{b}}] + [C^{*,\tilde{b}} - C^{*,\ell}].$$

By Lemma II.13, $D^{*,\bar{b}} - C^{*,\bar{b}} \leq 0$. The upper bound in Proposition II.12 then immediately follows from Lemma II.17 and Lemma II.16.

As for the lower bound for $C^{*,\ell}$, note that $C^{*,\ell} \geq C^{*,\bar{b}} - (c + L \cdot p) \cdot T \cdot \alpha \cdot \bar{D} \geq D^{*,\bar{b}} - (c + L \cdot p) \cdot T \cdot \alpha \cdot \bar{D}$, where the first inequality is from Lemma II.16. By the constraints in (2.16), we know that $\sum_{t=1}^T x_{t-L} \geq \sum_{t=1}^T \mu_t + w_T(\alpha)$, $z_{t+1} \geq \gamma_t(\alpha)$ for $1 \leq t \leq L+1$ and $z_{t+1} \geq w_t(\alpha)$ for $L+2 \leq t \leq T$. So, we can bound:

$$\begin{aligned}
C^{*,\ell} &\geq c \cdot \left[\sum_{t=1}^T \mu_t + w_T(\alpha) \right] + h \cdot \left[\sum_{t=1}^{L+1} \gamma_t(\alpha) + \sum_{t=L+2}^T w_t(\alpha) \right] - (c + L \cdot p) \cdot T \cdot \alpha \cdot \bar{D} \\
&= c \cdot \left[\sum_{t=1}^T \mu_t + \beta_{T-L}^{L+1}(\alpha) \right] + h \cdot \left[\sum_{t=1}^{L+1} \beta_1^t(\alpha) + \sum_{t=L+2}^T \beta_{t-L}^{L+1}(\alpha) \right] \\
&\quad - c \cdot \left[\sum_{t=1}^{T-1} \theta_t(\alpha) \right] - h \cdot \left[\sum_{t=1}^T \sum_{s=1}^{t-1} \theta_s(\alpha) \right] - (c + L \cdot p) \cdot T \cdot \alpha \cdot \bar{D} \\
&\geq c \cdot \left[\sum_{t=1}^T \mu_t + \beta_{T-L}^{L+1}(\alpha) \right] + h \cdot \left[\sum_{t=1}^{L+1} \beta_1^t(\alpha) + \sum_{t=L+2}^T \beta_{t-L}^{L+1}(\alpha) \right] \\
&\quad - 2 \cdot (c + h) \cdot \phi(T, \alpha) - (c + L \cdot p) \cdot T \cdot \alpha \cdot \bar{D}.
\end{aligned}$$

This completes the proof of Proposition II.12.

2.5 Conclusion

In this paper, we showed that stochastic inventory models with lead times and sequential probabilistic service level constraints can be well-approximated by deterministic programs when the targeted service level is sufficiently high. To accomplish this, we proposed a simple order-up-to control whose parameters can be computed using the solution of a deterministic approximation of backorder system, and proved its asymptotic optimality. Our work contributes to the growing body of literature that use asymptotic analysis to analyze the performance of simple heuristic controls in complex

stochastic inventory systems.

This work leaves several interesting future research directions. First, we mainly focus on the deterministic linear program approximation in this work, which is not necessarily tight due to the nature. To address this, instead of using a deterministic linear program, one could possibly use a more refined stochastic program or adapt a re-solving LP technique. We will leave the investigation on this direction for future research. Second, in our work, we have assumed that lead time is deterministic. It will be interesting, and impactful, if it can be shown that stochastic inventory models with random lead time and sequential probabilistic service level constraints can also be well-approximated by deterministic programs. Third, in this paper we have primarily focused on the analysis of order-up-to control. It is curious to see whether other simple heuristic controls such as constant-order policy is also asymptotically optimal (in some sense) in the presence of sequential probabilistic service level constraints. Note that the convexity argument in Xin and Goldberg [89] cannot be directly applied here. Fourth, it will be interesting to extend our analysis to other important neighboring inventory systems such as the dual-sourcing setting.

CHAPTER III

Snob and Follower Effects in Luxury Retailing

3.1 Introduction

The value of the luxury goods market is estimated at around 1.2 trillion Euro globally in 2018 (1.5% of GDP). The luxury industry has also grown rapidly, especially over the last two decades – for example, the global sales of personal luxury goods, a segment of the luxury products, has grown from 73 billion Euro in 1994 to 260 billion Euros in 2018 (Bain and Company [11]), which is double of the inflation rate. Although the moral legitimacy of luxury products is sometimes debatable, the luxury industry is too relevant for researchers to ignore, especially as it is becoming one of the drivers of economic growth (Berghaus, Müller-Stewens, and Reinecke [12]).

Luxury products has drawn attention from academia since Veblen [80]. It differs from other ordinary goods in many fundamental ways. Commonly, five effects are linked to luxury consumptions (Leibenstein [56]; Vigneron and Johnson [83]): (a) The Veblen effect, where people make conspicuous consumption to signal their wealth status; (b) the snob effect, where people value exclusiveness, scarcity, or uniqueness of a product; (c) the bandwagon effect, where people purchase for conformity by either following the purchase made by others in their group or imitating the affluent lifestyle of people whom they want to be (Dittmar [33]); (d) the hedonism effect, where people get emotional

satisfaction by purchasing; (e) the perfectionism effect, where people value the superior product quality or product characteristics.

Among the above five effects, the first three suggest an interesting structure of inter-personal influences (externalities) of the luxury products, namely the existence of both negative and positive externalities. The first and second effects reflect the *negative* externalities of luxury products – a higher sales of the product means either a less-efficient signal of wealth status or a dampened uniqueness; Customers’ utility of the product decreases in its total sales. We refer to such negative externalities as *snob effect*. The third effect indicates *positive* externalities: customers’ willingness-to-pay increase as he observes or expects higher purchases made by other people, either in his group or in a “higher-perceived” group. We refer to such positive externalities as *follower effect*.

The coexistence of these two opposing effects, i.e., follower and snob effects, are considered as the core of the luxury industry. We can easily find articles discussing how luxury brands deal with, or sometimes struggle with each of these effects. On the one hand, brands are working on enhancing the follower effect – for example, the CEOs of the leading luxury groups, LVMH group and the former Gucci group, view a core mission of luxury as selling dreams (Harvard Business Review, October 2001; Fortune, 6 September 2007). In other words, they advertise the consumption of the snobs and make more followers dream about the products and purchase them. On the other hand, to maintain the scarcity or exclusiveness of the products, brands might also want to limit the number of followers. For example, Tiffany & Co. decided to limit the sales of its “cash-cow” silver products to the massive volume of teen followers, in order to decrease the negative externalities to other (snob) customers (The Wall Street Journal, January 2007).

While the literature on luxury products includes extensive studies on the reasons

and drivers of luxury consumption (see Gurzki and Woisetschlager [42] for a comprehensive review), very few papers study the optimal selling strategies. At the same time, the literature that studies externalities contains extensive analysis of either positive or negative externalities, separately,¹ but very few papers study the snob and follower effects jointly and investigate how these two effects influence the retailing strategies under various market and pricing structures (more details in the Literature Review Section). In this paper, we aim to fill these gaps. We study the joint effect of snob and follower and analyze the selling strategies in the following three settings: (1) the product-line and pricing strategy in a monopoly market, where the retailer offers vertically-differentiated products, (2) the pricing strategies in a duopoly market where each retailer offers a single product, and (3) the product bundling strategies.

The aspects we consider are very relevant because we can easily observe a volume of interesting and (perhaps) counter-intuitive practices. Using personal fashion products as an example, we discuss these practices with respect to two scenarios: (1) products within the same category and (2) products across different categories. For scenario (1), we have two observations in cases where either a single retailer or multiple retailers offer various products. For the former case, the product-line strategy is relevant. For example, several luxury brands have secondary lines that offer similar types of products as the mainline but at lower prices (e.g., Prada has one secondary line MiuMiu, Ralph Lauren has Lauren Ralph Lauren and Polo Ralph Lauren as the secondary lines, and Armani (Giorgio Armani) holds secondary lines including Emporio Armani, Armani Exchange, etc.). Note that the general fashion brands like Zara and H&M, in contrast, usually only have one brand line. For the latter case, the competitive market structure can play a role. Note that a wide range of luxury products, with various qualities

¹On positive externalities: see Liebenstein 1950, Becker 1991, Besen and Farrell 1994, Katz and Shapiro 1994, Candogan, Bimpikis, and Ozdaglar [27], etc. On negative externalities: see Bagwell and Bernheim, 1996, Naor 1969, Lippman and Stidham 1977, Mendelson and Whang 1990, Momot et al. 2019, Tereyağođlu and Veeraraghavan [77] etc.

and prices, are offered by different companies – for example, in terms of handbags, such a range can be from higher price Chanel bags to lower price Tory Burch bags. The product range of general handbags has similarities, e.g., from relatively higher price handbags (with a general brand name) to lower price supermarket handbags. For scenario (2), bundling across products becomes an option. While the total mixed bundling strategy, where products are both offered separately and in bundles with discounts, is widely used for general products, it is rare to see such a form of bundling in the luxury industry. Instead, bundling in the luxury industry is often taken different forms, e.g., partial mixed bundling, where relatively low-value products are available both separately and in bundles, while a high-value product is only offered in the bundle. One example is that, for those highly-exclusive handbags like the Hermès Birkin and Kelly handbags, in most cases, they can only be obtained with the purchase of other relatively low value products (Forbes, January 2016; Bloomberg, June 2015). Another example is that, for certain limited-edition luxury sport cars like McLaren Speedtail, customers need to “spend on several regular models and ... be invited to put [their] name down for a limited-edition hypercar” (Economist, September 2018). Although this partial mixed bundling exist in practice, its underlining motivation is in fact not obvious – if the goal of such a bundling strategy is simply to generate a higher revenue through the bundle, then it can be easily achieved by charging a higher price for the high-value product instead.

Motivated by above observations, we aim to understand the driving forces from these three aspects of firm decisions through the lens of externalities. For each aspect, we use a stylized model to capture the main trade-offs of the market and product structures.

In the first aspect, product-line strategy, we model the retailer as a monopoly who offers two products with exogenous qualities. The monopolist decides which product(s)

to offer and the corresponding price(s). We have two findings. First, when the externalities are not too strong, the two opposite externalities, ie., snob and follower effects, turn out to reinforce each other and reduce the cannibalization between products. The optimal policy, thus, may be to offer both products and segment the market. As an comparison, the optimal policy in the no externalities case is to offer only one product. Second, when externalities become very strong, the cost of maintaining a self-selected segmentation becomes too large, and the optimal policy becomes to offer only one product.

In the second aspect, pricing in competitive market, we assume that each retailer offers a single product with exogenous quality. We find that, in equilibrium, while the two externalities work in the same direction, their impacts on market segmentation can be different depending on the distributions of customers' valuation. When customers' values follow a *discrete* distribution, the two externalities both induce a pricing war where one firm sells its products to both types of customers while the other cannot sell its product even by setting the price to 0. When customers' values follow a *continuous* distribution, the two externalities work through a joint mechanism to induce a new form of segmentation, i.e., partial segmentation, where one firm exclusively sells its product to one type of customers while the other one sells its product to both types of customers. In the continuous distribution case, we discover a new mechanism. Its intuition is as follows (more detailed discussions can be found in Section 3.5). On the one hand, a stronger follower effect indicates that one of the firms has the incentive to sell its products to more followers, since the followers' willingness-to-pay increases with the follower effect; such higher sales thus further reinforce the follower effect. On the other hand, a stronger snob effect indicates that the value of product offered by this firm is "contaminated" in the eyes of snobs, which makes it strictly dominated by the other product. The other firm, who faces less competition from the first firm, could

thus focus on snobs and charge a higher price.

In the third aspect, bundling strategy, to focus on the interaction between bundling decisions and externalities, we model the retailer as a monopoly. In terms of the optimal bundling policy, we find that it depends on the difference between customers' quality sensitivities. When snobs are much more sensitive to quality than followers, the retailer is more likely to offer the high-value product only in the bundle (partial-mixed bundling). Such a policy reduces the total sales of the high-value product and weakens the snob effect, but allow him to charge a higher price to snobs. When snobs' sensitivity to quality is similar to that of followers, the retailer is now more inclined to offer pure bundling – it is not wise to lose followers on any product in order to gain the “exclusiveness” from snobs who do not have a significant higher willingness-to-pay in such a case. In terms of the impact of externalities, we find that both the positive and negative externalities make the difference in quality sensitivities to be smaller and thus, make the retailer favor the pure bundling strategy.

Note that the insights related to the effect of joint positive and negative externalities are not limited to the settings of luxury products. These insights would apply to any product which features both snob and follower effects (e.g., customers purchase choices can be affected by the word-of-mouth effect, while many customers may still value the uniqueness of a non-luxury product to show their taste), or any service where additional joining customer creates both positive and negative externalities (e.g., more customers waiting for service may signal high quality of the service, i.e., positive externalities, but also create more congestion in the system, i.e., negative externalities (Veeraraghavan and Debo [81]; Veeraraghavan and Debo [82])).

The structure of the rest of the paper is as follows. We briefly discuss the most-relevant literature in Section 3.2. In Section 3.3, we introduce the general model of products and customer utilities. We study the effect of externalities along the

three outlined questions and develop our main results in Sections 3.4, 3.5, and 3.6, respectively.

3.2 Literature Review

Our discussion of literature consists of three parts, each corresponding to one of the aspects of the market and product structure that we study in this paper. For brevity, in this section, we only include the most relevant papers that consider externalities.

The product-line decision aspect is studied in both monopolistic and competitive settings. In monopolistic setting, firm decides whether to offer multiple products, or not, but only a few papers study the effect of externalities. They either include only snob effect or only follower effect: Jing [51] considers only the follower effect and studies how the effect changes the product portfolio decision. While in Jing [51], the follower effect is modeled in an aggregated way such that the total sales of all products affect the utility of each products, in our paper, both of the two externalities are modeled as product specific – the utility of a product is only affected by its own sales. Such product specific externalities not only better fit the luxury product settings we consider, but also change the structure of the trade-offs in the selling decision, and, thus provides new insights. Balachander and Stock [10] focus on the snob effect and investigate its impact on the decision about whether to offer limited edition product or not.² While, in Balachander and Stock [10], all customers either strictly prefer the limited edition product or are all indifferent between the two products (since the two products (normal and limited-edition products) the firm can potentially provide are always with the same quality and they only differ in whether there is a pre-announced amount of sales), in our paper, we consider a more general setting where different types of customers can

²While Balachander and Stock [10] mainly discuss the competitive setting of horizontal-differentiation case, it also covers a monopoly setting where horizon differentiation is mathematically equivalent to vertical differentiation.

have various preference of products (since we assume that the two products can be of heterogeneous quality and the sales of each product is generate by customers' choice instead of the exogenous limited-edition sales number). The only paper we are aware of, that has both externalities in monopolistic setting is Amaldoss and Jain [7] who study the branding decision, i.e., whether to sell two products under the same brand, facing both types of customer, snobs and followers. However, they do not allow customers to freely choose which product to purchase; Instead, they assume that any type of customer (snobs or followers) can only purchase the product that is exclusively offered for his own type. Thus, the paper does not reflect the critical nature of cannibalization across products in the vertical differentiation setting. In this paper, we include both snobs and followers effect, and further, we allow customers to choose the product that has a higher utility for him.³

In the competitive setting, while there is a rich literature on the network effect of products, there are very few papers consider both positive and negative externalities (explicitly as functions of product sales) and vertical-differentiated products. Many early papers (e.g., Matutes and Regibeau [60] and Laffont, Rey and Tirole [54]) focus on the compatibility of products offered by different firms – although they study products that have positive network effect, they do not consider the customers' utilities of products as functions of total sales. Instead, compatibility only means that customers have a larger choice set to “mix and match.” Among the papers where product utilities are explicitly modeled as functions of total sales, many of them use Hotelling model (of horizontal differentiation product) to model the intrinsic value of products (e.g., Balachander and Stock [10], Viswanathan [84]). Note that in this paper, however, we focus on the vertical differentiation, which better fits the luxury industry that

³Other two papers in this context that include both snob and follower effects are Amaldoss and Jain [4] and Amaldoss and Jain [6], both of which consider a single product – The former focuses on the slope (upward or downward) of the demand function, while the latter considers pricings in two periods where the snobs and followers make purchase sequentially in separate periods.

we study because customers highly value qualities and there is usually a clear rank across different products (a key feature in vertical differentiation). The key difference between horizontal and vertical differentiation product is that the former implies negative correlations of customers' utility across products while the latter implies positive correlations. Thus, the insights generated from these two differentiation settings could be very different. So far we know, the most relevant papers are Katz and Shapiro [52], Balachander and Stock [10], and Amaldoss and Jain [5]. Katz and Shapiro [52] captures the symmetric equilibrium in a Cournot duopoly game with (explicit) positive externalities. While symmetric equilibrium can be representative for a single follower effect, asymmetric equilibrium, however, is inevitable for the asymmetric externalities (snobs and followers) we study in this paper (we discuss in more details later in Section 3.5). Balachander and Stock [10] study how competition influences the decision whether to offer limited-edition products or not. While Balachander and Stock [10] only consider negative externalities, we expect that including both positive and negative externalities can bring new insights on the competitive market (remember that one of the key lessons from the externalities literature is that the positive externalities may change the equilibrium). Amaldoss and Jain [5] considers both follower and snob effects and focuses on the horizontally-differentiated products. Similar to the discussion of horizontal differentiation above, we expect that the insights are different in the vertical differentiation setting we consider in this paper.

The bundling strategy is studied in Marketing, Economics, Information Systems, and Operations Management literatures. The economic intuition for bundling is based on extracting a higher surplus when consumers have heterogeneous valuations of products (Adams and Yellen [2], Schmalensee [69]). We refer readers to Stremersch and Tellis [73] for a comprehensive review of bundling. For the bundling strategy of vertically differentiated products, which we study in this paper, the most relevant paper is

Ma and Mallik [59], but they do not consider any externalities. Another related paper is Prasad, etc. [64], where they provide the optimal bundling strategy with positive externalities. In this paper, we focus on positive and negative externalities jointly when studying the optimal bundling strategy. Another main difference between our paper and Prasad, etc. [64] is that while Prasad, etc. [64] is mainly about the interactions of costs and positive externalities, we focus on the zero-cost cases and show that the two externalities themselves can already change the optimal bundling strategies.

Our work also contributes to the broader modeling literature of luxury products. While early modeling papers focus on the drivers of the consumption of luxury products (Corneo and Jeanne [31]) and the reasons why brands leverage on the product scarcity (Stock and Balachander [72]), recent works have been shifting their focus to the optimal decisions and strategies. For example, Rao and Schaefer [66] characterize the product depreciation and upgrade decision over time for a single product that customers use to signal their wealth status and, Momot, Belavina, and Girotra [61] characterize the decision of choosing which customer(s) to offer the product given their network structures, in the case where all customers value the exclusiveness. Our paper complements these recent studies by providing optimal decisions facing three different market and product structures mentioned above.

3.3 Model

In this section, we introduce the general model.

We consider two products that are differentiated in their qualities, $q_1 \geq q_2$, where qualities are given exogenously.⁴ Such a setting can be viewed as the second stage of decisions after firms have decided the quality. We focus on the second stage to study the effect of externalities, without interactions with costs in the stage of designing

⁴The assumption follows Balachander and Stock [10] and Shaked and Sutton [70]

product quality.

The customer utility of a product comes from two parts, the functional utility determined by the attributes of the product, and the network utility determined by only externalities. Depending on whether customers experience positive or negative externalities, we consider two types of customers: snobs and followers. Snob customers are those who experience negative externalities – they value the exclusiveness and uniqueness of products. Thus, their utility decreases with total sales. Follower customers experience positive externalities – they like to follow the purchase of others and, thus, their utility increases with total sales. We normalize the quantity of snobs to 1 unit and use β to denote the quantity of followers. The utilities of snobs and followers are assumed to be:

$$U_i^S(v) = vq_i - kD_i^e - p_i$$

$$U_i^F(v) = vq_i + mD_i^e - p_i$$

where p_i is the price of product i and D_i^e represents the expected sales of product i .

For a type- v customer, where v represents his sensitivity to quality, the functional value of product i is vq_i . For snobs, the network value is $-kD_i^e$, where k represents the sensitivity to uniqueness. A higher total sales of product i hurts the product uniqueness and, thus, results in a lower utility for a snob customer. Similarly, for a follower customer the network utility is mD_i^e , where m represents the sensitivity to conformity. A higher total sales of product i increases his utility.

Note that we express the network value of a product as a function of total sales. This allows to capture the main drivers of positive and negative externalities (snobs and followers) and avoids the complexity of considering the structure of the social network. Such an approach is also one of the most widely-used ones in the literature,

e.g., Prasad, etc. [64], Amaldoss and Jain [7].

3.4 Product-line Strategy

In this section, we consider a simple case, where a monopolist determines the prices of two products with exogenous qualities, and illustrate how the snob and follower effects affect the pricing and the product-line decisions.

We assume that each customer needs only one unit of the product, and purchases the product with higher utility. Customers have an outside option, which we normalize to zero, without loss of generality.

We first consider the case, where all snobs are of the same type v_S , with mass unit one, and all followers are of the same type v_F , with mass β . Such discrete distribution setting is widely used in the literature (e.g. Adams and Yellen [2]; Stremersch and Tellis [73]), especially in the stage of illustrating the main trade-offs. We adopt this setting to illustrate the main changes brought by both positive and negative externalities. We assume that $v_S > v_F$, i.e., snobs have higher willingness-to-pay (Vigneron and Johnson [83]).

We start the analysis with a benchmark case where there are no externalities.

No Externalities

With no externalities, the utilities of customers are reduced to the following form.

$$U_i^S(v_S) = v_S q_i - p_i$$

$$U_i^F(v_F) = v_F q_i - p_i$$

Although there is no snob or follower effect in this case, we still denote the customers of type v_S as snobs and those of type v_F as followers, to be consistent with the rest of the paper.

The decision rule for each type of customer can be easily calculated. Snob chooses to purchase product 1 (the higher-quality product) if $v_S \geq \frac{p_1 - p_2}{q_1 - q_2}$ and $v_S \geq \frac{p_1}{q_1}$ (such that $v_S q_1 - p_1 \geq v_S q_2 - p_2$ and $v_S q_1 - p_1 \geq 0$), while he chooses to purchase product 2 if $v_S \leq \frac{p_1 - p_2}{q_1 - q_2}$ and $v_S \geq \frac{p_2}{q_2}$. Similarly, for a follower, he purchases product 1 if $v_F \geq \frac{p_1 - p_2}{q_1 - q_2}$ and $v_F \geq \frac{p_1}{q_1}$ or product 2 if $v_F \leq \frac{p_1 - p_2}{q_1 - q_2}$ and $v_F \geq \frac{p_2}{q_2}$. With this decision rule, there does not exist a market outcome such that high quality sensitive customers, i.e., snobs, purchase the high-quality product (product 1) while low quality sensitive customers (followers) purchase the low-quality product (“market segmentation”).

Proposition III.1. When there are no externalities, under the optimal prices, the sales outcome can be either of the two following cases:

- If $v_S \geq (1 + \beta)v_F$, snobs purchase product 1 while followers make no purchase. The optimal prices are $p_1^* = v_S q_1$ and any $p_2^* > v_F q_2$.
- Otherwise, both snobs and followers purchase product 1. The optimal prices are $p_1 = v_F q_1$ and any $p_2 > v_F q_2$.

In other words, although there are two potential products with different qualities, under the optimal policy, the firm only sells one product. Such a policy induces purchase from either both snobs and followers or only snobs, depending on customers’ sensitivity to quality.

The intuition is that, if a firm wants to offer both products, these two products naturally cannibalize each other. In order to prevent snob customers from purchasing the low-quality product (instead of the high-quality one), the price of the high-quality product needs to be set lower than its actual value to snobs. If the followers’ willingness-to-pay for the low product is low enough ($v_S \geq (1 + \beta)v_F$), then the cannibalization effect dominates the revenue from the low-quality product. Thus, it is optimal to only attract snobs customers to purchase the high-quality product. On the other hand, for

the case where the followers are sensitive enough to quality ($v_S < (1 + \beta)v_F$), if a firm is better off by lowering prices and bear cannibalization to allow the followers to purchase the low-quality product (instead of excluding them from any product), then the firm can be even better off by lowering prices all the way to allow followers purchasing the high-quality product. This is due to the multiplicity structure of the functional value, i.e., $vq_i \forall i \in \{1, 2\}$.

The analysis in Proposition III.1 can be extended to the case where multiple products are offered and multiple choices of customers.

In the rest of the paper, we denote $\Delta q = q_1 - q_2$.

Snob and Follower Effects

With both snob and follower effects, the utility of one type of customer is affected by the purchase of the other type of customer.

For the snob customers to purchase product 1, the utility of product 1 needs to be positive to them, $v_S q_1 - kD_1^e - p_1 \geq 0$, and also be larger than that of product 2, $v_S q_1 - kD_1^e - p_1 \geq v_S q_2 - kD_2^e - p_2$:

$$\left\{ \begin{array}{l} v_S \geq \frac{p_1 - p_2 + k(D_1^e - D_2^e)}{\Delta q} \\ v_S \geq \frac{p_1 + kD_1^e}{q_1} \end{array} \right.$$

For the snob customers to purchase product 2, we need $v_S q_2 - kD_2^e - p_2 \geq 0$ and $v_S q_1 - kD_1^e - p_1 < v_S q_2 - kD_2^e - p_2$:

$$\left\{ \begin{array}{l} v_S \geq \frac{p_2 + kD_2^e}{q_2} \\ v_S < \frac{p_1 - p_2 + k(D_1^e - D_2^e)}{\Delta q} \end{array} \right.$$

Similarly, the followers choose to purchase product 1 if

$$\begin{cases} v_S & \geq \frac{p_1 - p_2 - m(D_1^e - D_2^e)}{\Delta q} \\ v_S & \geq \frac{p_1 - mD_1^e}{q_1} \end{cases}$$

And the followers purchase product 2 if

$$\begin{cases} v_S & < \frac{p_1 - p_2 - m(D_1^e - D_2^e)}{\Delta q} \\ v_S & \geq \frac{p_2 - mD_2^e}{q_2} \end{cases}$$

With both externalities, we find that a market segmentation becomes possible, which is in contrast with the “no segmentation” result with no externalities.

Theorem III.2. *When there are both positive and negative externalities, under the optimal prices, the equilibrium can be of the following forms depending on the values of parameters:*

(1) *if $(v_S - v_F)q_2 \leq (m + k)\beta$, the equilibrium can be in the following four forms*

- *Snob customers purchase product 1 while follower customers purchase product 2 if the condition $(v_S - v_F)(\Delta q) \geq (k + m)(1 - \beta)$ holds. The optimal prices are $p_1 = v_S q_1 - k$ and $p_2 = \min\{v_F q_2 + m\beta, v_S q_2 - v_F(\Delta q) - k - m(1 - \beta)\}$ and the revenue is $R_1 = p_1 + \beta p_2 = v_S q_1 - k + \beta \min\{v_F q_2 + m\beta, v_S q_2 - v_F(\Delta q) - k - m(1 - \beta)\}$.*
- *Both snob and follower customers purchase product 1. The optimal price of product 1 is $p_1 = \min\{v_S q_1 - k(1 + \beta), v_F q_1 + m(1 + \beta)\}$. The price of product 2 can be set as long as $p_2 \geq p_1^* - \min\{v_S(\Delta q) - k(1 + \beta), v_F(\Delta q) + m(1 + \beta)\}$ is satisfied. The revenue is $R_2 = (1 + \beta) \min\{v_S q_1 - k(1 + \beta), v_F q_1 + m(1 + \beta)\}$*
- *Snob customers purchase 2 while followers purchase product 1 if the condition $(v_S - v_F)(\Delta q) < (k + m)(\beta - 1)$ holds. The optimal prices are $p_2^* = v_S q_2 - k$*

and $p_1^* = \min\{v_S q_2 - k + v_F \Delta q + m(\beta - 1), v_F q_1 + m\beta\}$. The optimal revenue is $R^* = v_S q_2 - k + \beta \min\{v_S q_2 - k + v_F \Delta q + m(\beta - 1), v_F q_1 + m\beta\}$.

- *Snob customers make no purchase while follower customers purchase product 1 if the condition of $(v_s - v_F)q_1 < (k + m)\beta$ also holds. The optimal prices are $p_1 = v_F q_1 + m\beta$ and any $p_2 \geq v_F q_2$. And $R_3 = \beta(v_F q_1 + m\beta)$*

(2) *If $(v_S - v_F)q_2 > (m + k)\beta$, the equilibrium can be in the following three forms*

- *Snob customers purchase product 1 while follower customers purchase product 2. The optimal prices are $p_1^* = v_F q_2 + m\beta + v_S \Delta q - k(1 - \beta)$ and $p_2^* = v_F q_2 + m\beta$. The revenue is $R^* = v_F q_2 + m\beta + v_S \Delta q - k(1 - \beta) + \beta(v_F q_2 + m\beta)$.*
- *Both snob and follower customers purchase product 1. The optimal prices and revenue are the same as shown above.*
- *Snob customers purchase 1 while follower customers make no purchase if the condition of $(v_S - v_F)q_1 \geq k + m$ holds. The optimal prices are $p_1^* = v_S q_1 - k$ and any $p_2 \geq v_S q_2$. The revenue is $R = v_S q_1 - k$.*
- *Snob customers purchase 2 while followers purchase product 1. The optimal prices and revenue are the same as shown above.*

The actual equilibrium outcome would fall into only one of the above scenarios depending on the parameter values.

Note that the equilibrium structure shown in Theorem III.2 is different from that shown in Proposition III.1 (the no externalities case). First, the market can be segmented where snobs purchase the high-quality product and followers purchase the low-quality one.⁵ Second, when the impact of the network values dominates the impact of

⁵Mathematically, it is because the externalities reverse the order of preference of the low-quality product. $v_S q_2 - \beta \leq v_F q_2 + m\beta$. Such a case can appear when either type difference is low, or q_2 is low, or any one of the snob or follower effect is large.

the functional values for all products⁶, equilibrium outcomes such that “only followers make purchase” or “followers purchase the high-quality product while snobs purchase the low quality one” can exist. To make a fair comparison to the no externalities case, we do not focus on such cases in the following analysis. In the rest of the paper, we make the following assumption to focus only on the cases where the externalities do not completely reverse the order of preference.

Assumption III.3. $v_S q_1 - k(1 + \beta) \geq v_F q_1 + m(1 + \beta)$.

Under Assumption III.3, the case (1) in Theorem III.2 is reduced to the following form.

Proposition III.4. Under Assumption III.3, the equilibrium in the case of $(v_S - v_F)q_2 \leq (m + k)\beta$ can be of the following forms.

- Snobs purchase product 1 while followers purchase product 2. The optimal prices are $p_1 = v_S q_1 - k$ and $p_2 = \min\{v_F q_2 + m\beta, v_S q_2 - v_F \Delta q - k - m(1 - \beta)\}$ and the revenue is $R = v_S q_1 - k + \beta \min\{v_F q_2 + m\beta, v_S q_2 - v_F \Delta q - k - m(1 - \beta)\}$.
- Both snobs and followers purchase product 1. The optimal prices are $p_1 = v_F q_1 + m(1 + \beta)$ and any $p_2 \geq p_1^* - \min\{v_S \Delta q - k(1 + \beta), v_F \Delta q + m(1 + \beta)\}$. The revenue is $R = (1 + \beta)[v_F q_1 + m(1 + \beta)]$

It is also easy to see that when $m = k = 0$, the equilibrium structure reduces to the no-segmentation form in Proposition III.2.⁷

We have two observations about how the externalities affect the equilibrium structure. First, when m or k increases, the optimal product strategies first change from

⁶In this scenario, the “willingness-to-pay” of snobs becomes smaller than that of the followers for all products, i.e., $v_S q_1 - k\beta < v_F q_1 + m\beta$.

⁷This is because that the feasibility condition of the last form of equilibrium does not hold. The remaining forms of equilibrium are the same as those in Theorem III.2.

the no-segmentation form (similar to Proposition III.1 without externalities) to “offering both product 1 and 2” (the first equilibrium shown in Proposition III.4), and then to “only offering product 1” (the second equilibrium shown in Proposition III.4). The intuition is that when m or k increases, it first eases the cannibalization of the two products, which makes a market segmentation profitable. However, when m and k becomes larger, followers’ utility of product 1 becomes closer to that of snobs and therefore, the “cost” of maintaining a market segmentation, i.e., keeping followers away from product 1, becomes higher. Note that although the positive and negative externalities affect the product utility in opposite directions, in the settings that we consider, they actually work in the same way by pushing the difference between utilities of snobs and utilities of followers smaller. Another observation is that when the portion of followers, β , increases, the optimal policy is more inclined to be “only offering product 1.” The intuition is similar to above; when the portion of followers increases, the profit losses from the cannibalization for followers become larger. This drives the firm to break the market segmentation and induce both snobs and followers to purchase product 1.

Our finding shares some similarities with Jing [51], where they also show that positive externalities can favor the market segmentation and change the product line offered by a monopoly. Our work complements their findings by further including the negative externalities in the analysis. We also enrich their findings: first, we derive additional insights into the product strategy changes by introducing the heterogeneity of quality sensitively across the customers; second, we show that both positive and negative externalities favor the market segmentation in the same direction, which may be counterintuitive at first glance (as one may suspect that negative externalities may go against market segmentation after knowing the positive-externalities results in Jing [51]).

3.5 Selling Strategy in Competitive Market

In this section, we focus on a competitive market where a single type of products are offered by different brands (with different qualities). Our interests in such a market structure stem not only from the industry practice as discussed in Section 3.1 but also from the insights we learn from Section 3.4 that externalities reduce the cannibalization between products. We are interested in whether externalities change the product decision in the (competition) settings where cannibalization from the same seller is irrelevant.

We study the duopoly market where each retailer $i \in \{1, 2\}$ offers a single product, with an exogenous quality, q_1 and q_2 ($q_1 \geq q_2$). Each retailer decides the price p_i for its product i . Without loss of generality, we assume that if a customer's utilities of the two products happen to be the same, he would always choose product 1.

In order to illustrate the impact of externalities, similar to Section 3.4, we start from studying the equilibrium without externalities as the benchmark, and then investigate the cases with both snob and follower effects. Note that the game structure is different in these two cases. The former can be directly analyzed using the standard Nash Equilibrium framework, while the latter needs to be captured by a joint framework of Rational Expectation and Nash Equilibrium. The intuition is that customers' expected demand and their corresponding purchase decisions would affect firms' pricing decisions and the equilibrium, and at the same time, the resulting equilibrium would affect whether the realized demands are consistent with customers' expectations (more detailed discussions on game structure will be provided in each subsection.)

In each of the no-externalities cases and with-externalities cases, we consider two types of customer distributions: discretely distributed customers, similar to that in Section 3.4, and continuously distributed customers. The reason that we add the latter type is as follows. In the competitive market that we study in this section, each

retailer i effectively faces only two possible scenarios (i.e., either snobs or followers choose to purchase product i), while, in the monopoly case in Section 3.4, the retailer faces four possible scenarios (i.e., snobs can purchase products 1 or 2 and followers can purchase products 1 or 2 as well). To make our results more robust, we thus further study uniformly distributed customers in addition to discretely distributed customers. These two distributions can be viewed as two extreme cases of a realistic customers' value distribution to help us draw insights of the general case. We also illustrate the new insights brought by continuously distributed customers by comparing the discrete and continuous cases.

The structure of this section is as follows. We discuss 4 different parts in total, which are combinations of the cases with and without externalities, and the cases with discrete and continuous customer distributions. In each part, we start with describing customers' decision rule (of which product to purchase) and the best response functions of both firms. We then characterize the equilibrium and investigate the impact of externalities and other relevant parameters.

3.5.1 No Externalities

We study the equilibrium with no externalities using the Nash equilibrium framework, focusing on the competition between two firms.

Discrete Distribution Cases

We start with the discrete distribution case. Customers' purchase decisions are as follows. For a customer with type v , either snobs or followers, he purchases product 1 from firm 1 if the following conditions hold: (1) $v \geq \frac{p_1 - p_2}{\Delta q}$ (such that $q_1 v - p_1 \geq q_2 v - p_2$) and (2) $v \geq \frac{p_1}{q_1}$. Similarly, he purchases product 2 from firm 2 if the following conditions hold: (1) $v < \frac{p_1 - p_2}{\Delta q}$ and (2) $v \geq \frac{p_2}{q_2}$.

Accordingly, we can describe the best response functions of the firms as follows:

(1) Given p_2 , the revenue of firm 1 is,

$$\pi_1 = \begin{cases} p_1 & \text{if } \min\{v_F(q_1 - q_2) + p_2, v_F q_1\} \leq p_1 \leq \min\{v_S(q_1 - q_2) + p_2, v_S q_1\} \\ (1 + \beta)p_1 & \text{if } p_1 \leq \min\{v_F(q_1 - q_2) + p_2, v_F q_1\} \end{cases}$$

(2) Given p_1 , the revenue of firms 2,

If the condition $\min\{p_1 - v_F(q_1 - q_2), v_F q_2\} \leq \min\{p_1 - v_S(q_1 - q_2), v_S q_2\}$ holds, then

$$\pi_1 = \begin{cases} p_2 & \text{if } \min\{p_1 - v_F(q_1 - q_2), v_F q_2\} \leq p_2 \leq \min\{p_1 - v_S(q_1 - q_2), v_S q_2\} \\ (1 + \beta)p_2 & \text{if } p_2 \leq \min\{p_1 - v_F(q_1 - q_2), v_F q_2\} \end{cases}$$

Otherwise, $\pi_1 = (1 + \beta)p_2$.

Proposition III.5. When firm 1 and 2 compete on price, given their products with quality q_1 and q_2 respectively, the equilibrium can be in either of the following forms

- If $(v_S - v_F)\Delta q \geq v_F(\beta q_1 - (1 + \beta)q_2)$, then firm 1 set price $p_1 = v_S\Delta q + v_F q_2$ and effectively only sell to snob customers; firm 2 set price $p_2 = v_F q_2$ and effectively only sell to follower customers.
- Otherwise, firm 1 set price $p_1 = v_F\Delta q$ and clear the market by selling to both snob and follower customers. Firm 2 generate no sales even by setting $p_2 = 0$.

Proposition III.5 shows that only when either snobs and followers have similar sensitivities to the quality or the two products have similar quality themselves, it is likely for firm 1 to fully occupy the market. The intuition is straightforward – when the customers and products become homogeneous, then the high-quality product strictly dominates the other and takes the full market.

For the first form of equilibrium in Proposition III.5, by comparing with Proposition III.1, we observe that competition brings back the perfect market segmentation – while

in the monopoly case the retailer does not induce snobs and followers to purchase different products, in the duopoly case each of the two retailers/competitors can sell the product exclusively to one type of customers. The intuition can be explained through two channels. The first channel is cannibalization: while cannibalization between two products prevents the monopoly from offering more than one product (as shown in Proposition III.1), such cannibalization within one retailer naturally disappears when the two products are managed by separate retailers. In other words, in the monopoly setting, providing an additional product cannibalizes with existing product(s) offered by the monopoly itself, while in the competitive setting, an additional product only cannibalizes with products offered by other firms. The second channel is through coordination. In the monopoly setting, the seller can coordinate two products and choose whether to offer only one of them or both. In the competitive setting, however, such coordination is not possible (and we do not consider any coordination game). These two channels lead to market segmentation in the competitive setting. Also, note that the equilibrium prices under market segmentation are the same as those in Proposition III.1 (although such segmentation scenario is dominated in the case of Proposition III.1).

The second form of equilibrium in Proposition III.5 describes the situation where firm 1 and 2 are involved in a *pricing war*, in which case the price of each product is driven down by each other until p_2 equals zero. Note that it is possible for one firm to take the entire market because customers are discretely distributed.

In the rest of the paper, we simply refer to these two forms of equilibrium in Proposition III.5 as perfect segmentation and pricing war, respectively.

Continuous Distribution Cases

Next, we analyze the equilibrium where customers' quality sensitivities/valuations follow uniform distributions. Motivated by the fact that snobs normally have a higher

sensitivity to quality, we assume that $v_F \sim U[0, 1]$ and $v_S \sim U[0, M]$, where $M \geq 1$. For any customer with type v , his purchase decision is the same as that described in the discrete distribution case above (where we provide a general description of customers' decision rules).

The best response functions of the firms are as follows:

(1) Given p_2 , the revenue of firm 1 is

$$\pi_1 = \begin{cases} \max_{p_1} p_1 [(\beta + 1) - (\beta + \frac{1}{M}) \frac{p_1 - p_2}{\Delta q}] & \text{if } p_1 > \frac{p_2 q_1}{q_2} \\ \max_{p_1} p_1 [(\beta + 1) - (\beta + \frac{1}{M}) \frac{p_1}{q_1}] & \text{if } p_1 \leq \frac{p_2 q_1}{q_2} \end{cases}$$

(2) Given p_1 , the revenue of firms 2 is

$$\pi_2 = \begin{cases} \max_{p_2} p_2 (\beta + \frac{1}{M}) [\frac{p_1 - p_2}{\Delta q} - \frac{p_2}{q_2}] & \text{if } p_2 < \frac{p_1 q_2}{q_1} \\ 0 & \text{o.w.} \end{cases}$$

Before we further analyze the equilibrium structure, it is useful to first discuss the following property.

Lemma III.6. *When customers follow continuous distributions, no single firm can fully occupy the market and earn positive revenue.*

Lemma III.6 means that a pricing war, similar to that in Proposition III.5, cannot survive in an equilibrium. In other words, there does not exist any equilibrium where the firm who sells a lower quality product, i.e., firm 2, cannot generate positive sales even by setting $p_2 = 0$. The intuition is that in the case where all customers purchase product 1, firm 2 can always increase p_2 a bit, attract a small portion of customers, and thus generate a positive revenue. We show the detailed proof below.

Proof. Proof of Lemma III.6 We first show that for any given $p_1 > 0$, a pricing war cannot survive in an equilibrium. We prove it by contradiction. Suppose there exist

a case where for $p_2 = 0$, firm 1 still occupies the whole market with some strictly positive price $p_1 > 0$. In this case, for any type t ($t \geq 0$) of customers, he has a positive utility from product 1, $tq_1 - p_1 \geq 0$. In such a case, firm 2 can increase p_2 to $p_\epsilon (> 0)$ to attract customers of type v where $0 \leq v \leq \frac{p_1 - p_\epsilon}{\Delta q}$. Any customer with type v prefers product 2 over product 1 and extracts a positive utility from product 2, since $vq_2 - p_\epsilon \geq vq_1 - p_1 \geq 0$ (where the last inequality follows from the condition $tq_1 - p_1 \geq 0$ for any $t \geq 0$). Thus, firm 2 can always deviate from $p_2 = 0$ and get a positive revenue. Next, it is easy to see that if firm 1 sets $p_1 = 0$, he can take the entire market but only with zero revenue. Combining these two parts, we complete the proof. \square

Next, we describe the structure of the equilibrium,

Proposition III.7. With products quality q_1 and q_2 and two types of customers with quality sensitivity $v_F \sim U[0, 1]$ and $v_S \sim U[0, M]$, when firm 1 and 2 compete on price, in equilibrium, firm 1 sets $p_1 = \frac{\beta+1}{\beta+\frac{1}{M}} \frac{2q_1}{4q_1-q_2} \Delta q$ and firm 2 sets $p_2 = \frac{\beta+1}{\beta+\frac{1}{M}} \frac{q_2}{4q_1-q_2} \Delta q$.

Note that in equilibrium, both snobs and followers make positive purchases from both firms. In the rest of the paper, we refer to this form of market segmentation as *Mixed Segmentation*.

Two properties of the equilibrium structure are worth noting. First, the intuition about the existence of mixed segmentation in the equilibrium is as follows. The fact that customers are continuously distributed suggests that firms no longer need to face the binary choice between a pricing war and a perfect segmentation as they face in the discretely distributed customers case (see Proposition III.5). Second, regarding the optimal prices, product 1 has a higher unit price of quality, $\frac{p_1}{q_1} \geq \frac{p_2}{q_2}$ (specifically, the relation of optimal prices in the equilibrium follows $\frac{q_1/p_1}{q_2/p_2} = \frac{1}{2}$). In other words, to achieve an additional unit of quality, customers need to pay more for product 1 relative to product 2. The intuition is that two firms divide the market share through different

prices per quality – firm 2 targets at those customers with low sensitivity by charging a lower price per unit quality while firm 1 targets at those high sensitivity customers. Such intuition is consistent with that in the market segmentation literature.

Another important takeaway is that no symmetric equilibrium exists even with symmetric qualities $q_1 = q_2$, or symmetric distributions where $M = 1$. The asymmetric structure of the equilibrium is inevitable given the asymmetric externalities (snobs and followers) we study in this paper.

In summary, while a perfect segmentation no longer exists in the equilibrium due to the nature of the continuity of customers, the equilibrium with continuously-distributed customers still has the same flavor of market segmentation as that in the discretely distributed customers case studied in Proposition III.5.

3.5.2 Snob and Follower Effects

With both positive and negative externalities, we need to adopt a joint framework of Rational Expectation and Nash Equilibrium in order to fully capture the game. This is because that customers' demand expectation and the resulting equilibrium have multiple layers of effects on each other – (1) firms' decisions on prices affect the demand expectations formed by customers, (2) customers' expected demand and the corresponding purchase decisions also affect firms' price choices and the resulting equilibrium, and (3) at the same time, the equilibrium also affects whether the realized demands is consistent with customers' expectations.

The sequence and structure of the game can be described as follows. First, each of the two firms sets the price individually. Then, after observing such prices, customers form rational expectations about the demands of both products and make purchase decisions correspondingly. Given such customers' purchasing behavior, each firm sets the price optimally according to its best response function, which is also coincident with

the price set at the beginning of the game. Then, in the equilibrium, importantly, the realized demand should also be consistent with the expectations formed by customers. Note that not every Nash Equilibrium satisfies the condition of equilibrium described above: if there is a Nash Equilibrium where the realized demands coincide with the demand expectations, we refer to it as a *rational equilibrium*; otherwise, even when a Nash equilibrium exists, we still consider it as a non-rational equilibrium because it violates the rational expectation condition.⁸

In the following analysis, we also follow Assumption III.3 in Section 3.4, such that the externalities do not completely reverse customers' product preferences.

Discrete Distribution Cases

We start with the discrete distribution case. Given any pair of prices (p_1, p_2) and the corresponding demand expectations (D_1^e, D_2^e) formed by customers, customers' choices between products are as follows:

For a snob customer with type v_S , he purchases product 1 from firm 1 if the following two conditions hold: (1) $v_S \geq \frac{kD_1^e + p_1}{q_1}$, and (2) $v_S \geq \frac{k(D_1^e - D_2^e) + p_1 - p_2}{\Delta q}$, meaning that the utility of product 1 is (1) positive (i.e., $q_1 v_S - kD_1^e - p_1 \geq 0$) and (2) higher than that of product 2 (i.e., $q_1 v_S - kD_1^e - p_1 \geq q_2 v_S - kD_2^e - p_2$.) Similarly, he purchases product 2 from firm 2 if $v_S \geq \frac{kD_2^e + p_2}{q_2}$ and $v_S \leq \frac{k(D_1^e - D_2^e) + p_1 - p_2}{\Delta q}$.

By the same token, for a follower customer with type v_F , he purchases product 1 from firm 1 if $v_F \geq \frac{-mD_1^e + p_1}{q_1}$ and $v_F \geq \frac{-m(D_1^e - D_2^e) + p_1 - p_2}{\Delta q}$. Similarly, he purchases product 2 from firm 2 if $v_F \geq \frac{-mD_2^e + p_2}{q_2}$ and $v_F \leq \frac{-m(D_1^e - D_2^e) + p_1 - p_2}{\Delta q}$.

The form of firms' best response functions vary across different customers' demand expectations.

Due to the discreteness of customer distribution, we only need to analyze best response functions with four different combinations of demand expectations: (1) $D_1^e = 1$

⁸Such a game structure is widely used in the literature of competition with externalities. See more discussions in Katz and Shapiro [52].

and $D_2^e = \beta$, where snobs purchase product 1 (from firm 1) while followers purchase product 2 (from firm 2), (2) $D_1^e = 1 + \beta$ and $D_2^e = 0$, where both snobs and followers purchase product 1, (3) $D_1^e = 0$ and $D_2^e = 1 + \beta$, where both snobs and followers purchase product 2, and (4) $D_1^e = \beta$ and $D_2^e = 1$, where snobs purchase product 2 while followers purchase product 1.⁹ We list out these four different scenarios to highlight the fact that not each Nash Equilibrium, with given demand expectations, would lead to a rational equilibrium – see discussion after Theorem III.8 below. The specific form of each firm’s best response, however, can be easily derived under each scenario. To reserve space, we relegate the full analysis to the proof of Theorem III.8 as shown in the Appendix.

Next, we describe the structure of the rational equilibrium.

Theorem III.8. *With both positive and negative externalities, there are two possible rational equilibria:*

- *If $(v_S - v_F)\Delta q \geq (m + k)(1 - \beta)$ and $\min\{v_S q_1 - k - (v_S - v_F)q_2 + \beta(k + m), v_S q_1 - k\} \geq (1 + \beta) \min\{v_S(q_1 - q_2) - k(1 + \beta), v_F(q_1 - q_2) + m(1 + \beta)\}$, firm 1 only sells to snobs while firm 2 only sells to followers. The prices set by firms are*

$$p_1 = \begin{cases} v_S q_1 - k & \text{if } (v_S - v_F)q_2 \leq \beta(m + k) \\ v_S q_1 - k - (v_S - v_F)q_2 + \beta(k + m) & \text{o.w.} \end{cases}$$

$$p_2 = v_F q_2 + m\beta$$

- *Otherwise, only firm 1 sells to both snobs and followers. The optimal price set by firm 1 is $p_1 = \min\{v_S(q_1 - q_2) - k(1 + \beta), v_F(q_1 - q_2) + m(1 + \beta)\}$. Firm 2 generates no sales even if setting $p_2 = 0$.*

⁹Note that there do not exist cases where only snobs (or followers) make a purchase (e.g., snobs purchase product 1 while followers make no purchase), as the firm of 0 sales can always attract the followers (or snobs) and be better off.

In terms of the equilibrium itself, we have two main observations. First, no rational equilibrium exists under demand expectations scenarios (3) and (4) described above (i.e., both snobs and followers purchase product 2 or snobs purchase product 2 while followers purchase product 1). In scenario (3), firm 1 always has the incentive to deviate from making zero sales (such that the realized demand of product 1 equals the expected demand, i.e., $D_1^e = 0$), while in scenario (4), given any p_2 , there is no such p_1 that equates the realized demand with the expected demand. Second, the structure of equilibrium with externalities (as shown in Theorem III.8) is the same as that without externalities (as shown in Proposition III.5).

While in the monopoly case externalities change the policy structure by reducing the cannibalization (as discussed in Section 3.4), they do not change the equilibrium structure in the competition setting where cannibalization or coordination is absent. In the competitive case, externalities, anyway, do not change the fact that firms need to make sharp choices between perfect segmentation and pricing war when facing discrete customers.

In terms of the impact of externalities and the portion of followers on equilibrium, we also have two main observations. The first observation is that, under most cases, both snob and follower effects are more likely to induce a pricing war in the equilibrium. We discuss the intuition in two parts. First, when k is away from a narrow medium-level range, i.e., $k \leq \frac{(v_S - v_F)\Delta q}{1 + \beta} - m$ or $k \geq \frac{(v_S - v_F)\Delta q}{1 - \beta} - m$, both the snob and follower effects make the seller with the higher-quality product (i.e., firm 1) have a higher incentive to induce purchases from all customers. This is because when either snob or follower effect is enhanced (k or m increases), followers' willingness-to-pays become closer to those of snobs. Therefore, it is less worthwhile for firm 1 to “lose” followers in exchange for snobs by creating “exclusiveness” values for them (note that snobs do not have a much higher willingness-to-pay than followers in this case). Second, when k is in a

narrow medium-level range, i.e., $\frac{(v_S - v_F)\Delta q}{1 + \beta} - m \leq k \leq \frac{(v_S - v_F)\Delta q}{1 - \beta} - m$, an increase in the snob effect, however, reduces the sellers' incentives to induce a pricing war. This is because k in this value range creates a distortion – price p_1 is more sensitive to snob effect in the pricing war compared with that in the market segmentation. If firm 1 wants to initiate a pricing war, the price that he needs to attract snobs is even lower than that to attract followers (remember that the “outside option” for any type of customers is to choose product 2 at no cost (i.e., $p_2 = 0$ in the pricing war)). Thus, an increase of snob effect means that it is more costly for firm 1 to initiate the pricing war by lowering p_1 .

The second observation is that when the portion of followers becomes larger, firm 1 has less incentive to ignore followers. Thus, it is more likely for firm 1 to go for the “price-war” equilibrium. The intuition is similar to that in Section 3.4.

Continuous Distribution Cases

Next, we analyze the equilibrium where customers' sensitivity to quality/valuations follow uniform distributions. Recall that we are only interested in any “rational” equilibrium where (1) the prices form a Nash Equilibrium and (2) realized demands coincide with expected ones.

Consistent with the continuous distribution cases in Section 3.5.1, we assume that $v_F \sim U[0, 1]$ and $v_S \sim U[0, M]$, where $M \geq 1$.

First, note that customers' decision rule is the same as that described in the discrete distribution case above (where we provide a general description of customers' decision rules) given prices (p_1, p_2) and their prior of demand (D_1^e, D_2^e) .

In order to capture the best response function of each firm, we first characterize the aggregate demand functions for both products for any given pair of prices (p_1, p_2) . There are two layers of complexity about demand functions. First, customers' demand may have various function forms with respect to different values of expected demand.

Second, although D_1^e and D_2^e are functions of p_1 and p_2 , there is no direct way for firms to pin down the exact demand expectations that customers form – Instead, both the expected and the realized demands need to be solved in a rational expectation equilibrium. Thus, similar to the discrete distribution case, we divide the analysis into 6 cases, depending on the given prices (p_1, p_2) and the possible demand expectations (see Table 3.1). It is worth noting that the cut-offs between different cases are not exogenous; instead, they are functions of choice variables p_1 and p_2 . Therefore, without solving for the optimal prices, firms do not know which cases it would fall into. This is one of the reasons why the equilibrium is challenge to capture and why we first need to analyze each case separately.

Case	Conditions	
1	$D_1^e q_2 - D_2^e q_1 \geq 0$	$p_2 q_1 - p_1 q_2 \geq k(D_1^e q_2 - D_2^e q_1)$
2		$-m(D_1^e q_2 - D_2^e q_1) \leq p_2 q_1 - p_1 q_2 \leq k(D_1^e q_2 - D_2^e q_1)$
3		$p_2 q_1 - p_1 q_2 \leq -m(D_1^e q_2 - D_2^e q_1)$
4	$D_1^e q_2 - D_2^e q_1 < 0$	$p_2 q_1 - p_1 q_2 \geq -m(D_1^e q_2 - D_2^e q_1)$
5		$k(D_1^e q_2 - D_2^e q_1) \leq p_2 q_1 - p_1 q_2 \leq -m(D_1^e q_2 - D_2^e q_1)$
6		$p_2 q_1 - p_1 q_2 \leq k(D_1^e q_2 - D_2^e q_1)$

Table 3.1: 6 Cases of Demand Functions

Before we get into the details of the aggregate demand functions, we first briefly discuss the meaning of two sets of conditions that define these six cases (i.e., two columns in Table 3.1). The conditions in the first column describe the relations between quality and demand. In the first 3 cases, for each unit of quality, firm 1 induces more demand than firm 2 does under the equilibrium, while in the last 3 cases, the reverse holds. The conditions in the second column describe the relation of price per quality between the two products. In cases 1 and 4, firm 1 always charges a lower price per quality than firm 2 does. In cases 2 and 5, either firm may charge a higher price per unit. In cases 3 and 6, firm 1 always charges a higher price per unit quality than firm 2 does.

The demand functions, D_1 and D_2 , in each of the above six cases are shown as below. Again, recall that as discussed above, we are only interested in the equilibrium where the expected demands coincide with the realized ones, i.e., $D_1 = D_1^e$ and $D_2 = D_2^e$. Also note that notation-wise, we still use both D_i and D_i^e ($i \in \{1, 2\}$) in the expressions below in case readers would like to see a clear link between the customers' decision rule and the aggregate demand functions.

Case 1: All customers, both snobs and followers, choose to purchase product 1, but not to purchase product 2.¹⁰

$$\begin{aligned} D_1 &= (\beta + 1) - \left(\beta + \frac{1}{M} \right) \frac{p_1}{q_1} + \left(\beta m - \frac{k}{M} \right) \frac{D_1^e}{q_1} \\ D_2 &= 0 \end{aligned} \quad (3.1)$$

Case 2: Followers only purchase product 1, while snobs purchase both products (some snobs choose product 1, while the others choose product 2).¹¹

$$\begin{aligned} D_1 &= (\beta + 1) - \beta \frac{-mD_1^e + p_1}{q_1} - \frac{1}{M} \frac{k(D_1^e - D_2^e) + p_1 - p_2}{\Delta q} \\ D_2 &= \frac{1}{M} \left(\frac{k(D_1^e - D_2^e) + p_1 - p_2}{\Delta q} - \frac{kD_2^e + p_2}{q_2} \right) \end{aligned} \quad (3.2)$$

Case 3: Both snobs and followers purchase both product 1 and 2 (for either snobs or followers, some of them choose product 1, while others choose product 2).¹²

¹⁰ $D_1^F = \beta(1 - \frac{-mD_1^e + p_1}{q_1})$, $D_2^F = 0$, $D_1^S = \frac{1}{M}(M - \frac{kD_1^e + p_1}{q_1})$, and $D_2^S = 0$.
¹¹ $D_1^S = \frac{1}{M}(M - \frac{k(D_1^e - D_2^e) + p_1 - p_2}{\Delta q})$, $D_2^S = \frac{1}{M}(\frac{k(D_1^e - D_2^e) + p_1 - p_2}{\Delta q} - \frac{kD_2^e + p_2}{q_2})$, $D_1^F = \beta(1 - \frac{-mD_1^e + p_1}{q_1})$, and $D_2^F = 0$.
¹² $D_1^F = \beta(1 - \frac{-m(D_1^e - D_2^e) + p_1 - p_2}{\Delta q})$, $D_2^F = \beta(\frac{-m(D_1^e - D_2^e) + p_1 - p_2}{\Delta q} - \frac{-mD_2^e + p_2}{q_2})$, $D_1^S = \frac{1}{M}(M - \frac{k(D_1^e - D_2^e) + p_1 - p_2}{\Delta q})$, and $D_2^S = \frac{1}{M}(\frac{k(D_1^e - D_2^e) + p_1 - p_2}{\Delta q} - \frac{kD_2^e + p_2}{q_2})$.

$$D_1 = (\beta + 1) - \beta \frac{-m(D_1^e - D_2^e) + p_1 - p_2}{\Delta q} - \frac{1}{M} \frac{k(D_1^e - D_2^e) + p_1 - p_2}{\Delta q} \quad (3.3)$$

$$D_2 = \beta \left(\frac{-m(D_1^e - D_2^e) + p_1 - p_2}{\Delta q} - \frac{-mD_2^e + p_2}{q_2} \right) + \frac{1}{M} \left(\frac{k(D_1^e - D_2^e) + p_1 - p_2}{\Delta q} - \frac{kD_2^e + p_2}{q_2} \right) \quad (3.4)$$

Case 4: All customers, both snobs and followers, choose to purchase product 1 but not to purchase product 2. Note that although the demand functions are the same as that in Case 1, the “feasible sets” are different – Case 1 and 4 are under different conditions of quality and expected demands.

Case 5: Snobs only purchase product 1, while followers purchase both product 1 and 2 (some followers choose product 1, while others choose product 2).¹³

$$D_1 = \beta + 1 - \beta \frac{-m(D_1^e - D_2^e) + p_1 - p_2}{\Delta q} - \frac{1}{M} \frac{kD_1^e + p_1}{q_1} \quad (3.5)$$

$$D_2 = \beta \left(\frac{-m(D_1^e - D_2^e) + p_1 - p_2}{\Delta q} - \frac{-mD_2^e + p_2}{q_2} \right)$$

Case 6: Both snobs and followers purchase both product 1 and 2. Note that although the demand functions are the same as in Case 3, the “feasible sets” are different – Case 3 and 6 are under different conditions of quality and expected demands.

Among these six cases, we first observe that certain demand functions are similar to the corresponding/benchmark cases as we discussed earlier. In particular, Case 1 and 4 are similar to the pricing war studied in Theorem III.8 (Discrete distribution case with externalities). Case 3 and 6 are similar to the mixed segmentation in Proposition III.7 (Continuous distribution case without externalities). Further, we also observe a

¹³ $D_1^F = \beta \left(1 - \frac{-m(D_1^e - D_2^e) + p_1 - p_2}{\Delta q} \right)$, $D_2^F = \beta \left(\frac{-m(D_1^e - D_2^e) + p_1 - p_2}{\Delta q} - \frac{-mD_2^e + p_2}{q_2} \right)$, $D_1^S = \frac{1}{M} \left(M - \frac{kD_1^e + p_1}{q_1} \right)$, and $D_2^S = 0$.

new form of segmentation that does not exist in either the no-externalities case or the discrete distribution case. That is, one type of customers exclusively purchase only one product, while the other type of customers purchase both products. There are two specific forms of such segmentation: (1) all followers only purchase product 1 (and snobs make positive purchases for both products), as described in case 2, and (2) all snobs only purchase 1 (and followers make positive purchases for both products), as described in case 5. We refer to such a segmentation as *Partial Segmentation* in the rest of the paper.

Based on the demand functions above, we can implicitly derive the following best response functions of firm 1 and 2. The complete form of the best responses can be found in the Appendix. Note that while the best response of each firm needs to be analyzed in each demand case separately, both firm 1 and 2 need to pick the same case in a rational equilibrium.

Before we characterize the equilibrium, we first show the following property.

Lemma III.9. *When customers' values follow uniform distribution, there does not exist any equilibrium where one firm fully occupied the market and generates positive revenue.*

Lemma III.9 means that any form of pricing war, as described in demand cases 1 and 4, could not survive in the rational equilibrium. The intuition is similar to that of Lemma III.6. We, however, use a different type of argument that fits the settings here. Consider the equilibrium in each of the two possible scenarios. In scenario 1, firm 1, knowing the fact that firm 2 will never set price as $p_2 \geq \frac{k(D_1^e q_2 - D_2^e q_1) + p_1 q_2}{q_1}$ (and thus get 0 revenue), will set price within the last two cases. In scenario 2, firm 1, knowing the fact that firm 2 will never set price as $p_2 \geq \frac{-m(D_1^e q_2 - D_2^e q_1) + p_1 q_2}{q_1}$ (and thus get 0 revenue), will also set price within the last two cases. We omit the full proof.

In terms of the equilibrium structure, we note that it is difficult to characterize it

analytically. As is discussed before, there are multiple layers of complexity: (1) Any equilibrium outcome needs to satisfy the conditions of both rational expectation and Nash Equilibrium, (2) there is no explicit demand function for any given prices (Instead, demand functions have various function forms depending on the expected demand, and we need to solve a rational expectation equilibrium to know the realized demand), and (3) both firm 1 and 2 need to pick prices that satisfy the same set of conditions (listed in Table 3.1) to form a rational equilibrium. Given such complexities, instead, we study the properties of equilibrium numerically. In the numerical experiments, the parameters we use are: $M \in \{1.5, 2\}$, $m \in \{0.1, 0.2\}$, $k \in \{0.1, 0.2\}$, $\beta \in \{0.5, 1.5\}$, and $q_2 \in \{0.3, 0.7\}$. Without loss of generality, we normalize q_1 to 1. The parameters are carefully chosen such that Assumption 1 is always satisfied – the preferences of products are not completely altered by the externalities.

Table 3.2 reports the equilibria. The last column list the corresponding realized demand structure, which refer to the six cases listed in Table 3.1.

Based on the simulation results above, we observe several properties of the equilibrium.

First, our main observation is that stronger snob and follower effects are more likely to induce partial segmentation in the equilibrium. We identify two mechanisms through which externalities can play a role. The first mechanism, formed jointed by snob and follower effects, was absent in the no-externalities or the discrete distribution cases, and therefore, brings new insight. We call it as the “competition-free” mechanism: On the one hand, when follower effect becomes stronger, firm 2 has more incentive to sell product 2 to followers since the followers’ willingness-to-pay increases with m and is further enhanced by a higher sales D_2 . On the other hand, when snob effect becomes stronger, the value of product 2 decreases for snobs (because of the higher sales of product 2) such that snobs prefer product 1 over product 2. Firm 1 could thus

case	Parameters						Equilibrium				Demand Case
	M	m	k	β	q_1	q_2	p_1	p_2	D_1	D_2	
1	1.50	0.10	0.10	0.50	1.00	0.30	0.50	0.08	0.79	0.39	6
2	1.50	0.10	0.10	0.50	1.00	0.70	0.25	0.09	0.87	0.47	3
3	1.50	0.10	0.10	1.50	1.00	0.30	0.38	0.07	1.66	0.50	3
4	1.50	0.10	0.10	1.50	1.00	0.70	0.37	0.19	1.53	0.42	2
5	1.50	0.10	0.30	0.50	1.00	0.30	0.55	0.10	0.67	0.30	6
6	1.50	0.10	0.30	0.50	1.00	0.70	0.34	0.15	0.68	0.47	6
7	1.50	0.10	0.30	1.50	1.00	0.30	0.44	0.07	1.32	0.58	6
8	1.50	0.10	0.30	1.50	1.00	0.70	0.45	0.29	0.44	0.41	2
9	1.50	0.30	0.10	0.50	1.00	0.30	0.43	0.07	1.00	0.30	3
10	1.50	0.30	0.10	0.50	1.00	0.70	0.38	0.18	0.89	0.35	2
11	1.50	0.30	0.10	1.50	1.00	0.30	0.69	0.40	0.81	1.69	5
12	1.50	0.30	0.10	1.50	1.00	0.70	0.70	0.54	0.50	0.97	5
13	1.50	0.30	0.30	0.50	1.00	0.30	0.51	0.08	0.76	0.36	6
14	1.50	0.30	0.30	0.50	1.00	0.70	0.46	0.27	0.70	0.33	2
15	1.50	0.30	0.30	1.50	1.00	0.30	0.72	0.40	0.64	1.52	5
16	1.50	0.30	0.30	1.50	1.00	0.70	0.51	0.47	0.55	1.41	5
17	2.00	0.10	0.10	0.50	1.00	0.30	0.57	0.08	0.81	0.40	6
18	2.00	0.10	0.10	0.50	1.00	0.70	0.28	0.11	0.41	0.47	3
19	2.00	0.10	0.10	1.50	1.00	0.30	0.42	0.07	1.68	0.50	3
20	2.00	0.10	0.10	1.50	1.00	0.70	0.39	0.21	1.64	0.35	2
21	2.00	0.10	0.30	0.50	1.00	0.30	0.62	0.10	0.71	0.33	6
22	2.00	0.10	0.30	0.50	1.00	0.70	0.36	0.16	0.74	0.46	3
23	2.00	0.10	0.30	1.50	1.00	0.30	0.43	0.06	1.46	0.62	6
24	2.00	0.10	0.30	1.50	1.00	0.70	0.45	0.27	0.44	0.38	2
25	2.00	0.30	0.10	0.50	1.00	0.30	0.47	0.07	1.04	0.31	3
26	2.00	0.30	0.10	0.50	1.00	0.70	0.42	0.19	0.94	0.33	2
27	2.00	0.30	0.10	1.50	1.00	0.30	0.74	0.40	0.87	1.63	5
28	2.00	0.30	0.10	1.50	1.00	0.70	0.61	0.49	0.66	1.26	5
29	2.00	0.30	0.30	0.50	1.00	0.30	0.57	0.08	0.81	0.40	6
30	2.00	0.30	0.30	0.50	1.00	0.70	0.49	0.27	0.33	0.32	2
31	2.00	0.30	0.30	1.50	1.00	0.30	0.71	0.40	0.83	1.67	5
32	2.00	0.30	0.30	1.50	1.00	0.70	0.55	0.47	0.63	1.37	5

Table 3.2: Equilibriums in Competitive Market

focus on snobs without worrying about the competition from product 2 and therefore, is able to charge a higher price p_1 . In all, stronger follower and snob effects are more likely to induce one form of partial segmentation, in which firm 2 only sells product 2 to followers while firm 1 sells product 1 to both snobs and followers (Demand Case 5), in the equilibrium. Such a mechanism plays a dominating role when the portion of followers is larger than the snobs and the quality between products is quite different (examples of such changes can be seen by comparing Case 11, 15 to Case 3, 7, and comparing Case 27, 31 to Case 19, 23, where $\beta = 1.5$ and $(q_1, q_2) = (1, 0.3)$).

The second mechanism follows a similar spirit as in Theorem III.8 (see the discussion of impacts of externalities on equilibrium). When follower and snob effects become stronger (i.e., m and k increases), the followers' willingness-to-pay for product 1 becomes closer to the snobs', and it drives firm 1 to sell product 1 to both types of customers. Given the high sales of product 1 and a large follower effect, followers naturally have little incentive to purchase product 2. As a result, another form of partial segmentation (in which firm 2 only sells product 2 to snobs while firm 1 sells product 1 to both snobs and followers) is thus more likely to occur in the equilibrium (Demand Case 2). Such a mechanism plays a dominating role when the portion of followers is smaller than the snobs and the quality of two products are similar (examples of such changes can be seen by comparing Case 10, 14 to Case 2, 6, and comparing Case 26, 30 to Case 18, 22, where $\beta = 0.5$ and $(q_1, q_2) = (1, 0.7)$).¹⁴

It is worth noting that when the two mechanisms above play dominating roles, a larger portion of followers further enhance the forces in these mechanisms and make it even more likely to result in a partial segmentation in equilibrium (for the form described in Demand Case 2, see examples by comparing Case 2 and 6 with Case 4 and

¹⁴Note that when β is small and the difference between q_1 and q_2 is small, there are not enough purchases from followers to boost the snob effect and not enough quality difference to construct a "competition-free" situation for firm 1. Thus, a segmentation similar to that of the first mechanism does not occur.

8, respectively; for the form described in Demand Case 5, see examples by comparing Case 10 and 14 to Case 12 and 16, respectively).

Some other impact of externalities includes (1) altering the relation of induced demand per quality between two products (when snob effect increases, the condition of demand changes from $D_1q_2 - D_2q_1 \geq 0$ to $D_1q_2 - D_2q_1 < 0$)¹⁵ and (2) affecting the optimal prices p_1 and p_2 (increase of k causes increases in both p_1 and p_2).¹⁶ We view these impacts as minor ones, as the structure of the equilibrium is unaffected, and provide brief discussions of intuitions in the footnotes above.

Another interesting observation is that, in general, the equilibrium structure is not affected by (1) the upper bound of the value distribution of snobs (i.e., M) and (2) the quality difference between two products (i.e., $q_1 - q_2$). This is because that the order of customers' willingness-to-pay distributions between snobs and followers and the preference of products are unchanged by the above parameters. Therefore, neither customers' decisions nor firms' decisions would be affected.

In sum, we identify two mechanisms where snob and follower effects work together to induce a partial segmentation in the equilibrium. While the portion of followers could also play a role, we find other parameters irrelevant in determining the equilibrium structure.

¹⁵When m is small, larger snob effect changes the equilibrium from case 3 to case 6, i.e., the condition of demand changes from $D_1q_2 - D_2q_1 \geq 0$ to $D_1q_2 - D_2q_1 < 0$ (Examples can be seen by comparing cases 6 and 7 to cases 2 and 3). Basically, when suffering more from the negative externalities (snob effect), firm 1 has less incentive to sell to more customers and thus sets a higher price p_1 , which reverses the relation of induced demand per unit quality – product 2 becomes more appealing to customers.

¹⁶When m is large and the equilibrium is already in the form of partial segmentation (as is shown in the main observation above), increase of k causes increases in both p_1 and p_2 (see examples by comparing case 14 and case 10 and comparing case 15 and case 11). The intuition is that a stronger snob effect makes firm 1 to induce fewer sales and sets a higher price. Firm 2, therefore, could also charge a higher p_2 without losing much market share.

3.6 Bundling Strategy

In the previous two sections, we consider the settings where the firm(s) sell one type (category) of product (either one firm sells two products under the same category in the monopoly setting, or two firms sell different products under the same category), and each customer only needs one unit of such a product. In this section, we consider the case where a firm offers two types (categories) of product (with exogenous quality q_1 and q_2),¹⁷ and each customer needs one unit of each type of product. We are interested in whether allowing the firm to offer bundles of products would affect its pricing and selling strategies.

The firm/retailer has three different types of selling strategies: (1) *Pure component*, where products are offered separately, with price p_1 and p_2 , respectively; (2) *Pure bundling*, where only a bundle containing both products is offered, with price p_B (no product is sold individually); (3) *Mixed bundling*, where both the bundle and product(s) are offered. Depending on the number of products offered individually, the mixed bundling has two forms. The first form is called partial mixed bundling, where one of the products is not offered individually. The second form is called total mixed bundling, where both of the products can be purchased individually.

We assume that customers' utilities from a bundle are additive, i.e., the value of a bundle is the sum of the value of each of its components. In other words, the two products are neither complements nor substitutes. Such a setting (where the only possible connection between the two products is through the bundle) allows us to concentrate on the effects of externalities on bundling strategy and avoid distortions from product substitution/compatibility. .

Regarding the network value, we assume that snobs are only affected by the snob

¹⁷The two-product setting is treated as a representative case to get insights for the more general cases (Stremersch and Tellis [73]).

effect for their purchase of product 1, i.e., $U_1^S(v_S) = v_S q_1 - kD_1^e - p_1$ and $U_2^S(v_S) = v_S q_2 - p_2$. Follower customers, on the contrary, still experience positive externalities on any of their purchases, i.e., $U_1^F(v_F) = v_F q_1 + mD_1^e - p_1$ and $U_2^F(v_F) = v_F q_2 + m_2^e - p_2$ (which is unlike that in the single product-category case where snobs experience negative externalities on any of their purchases). Note that such an asymmetric structure of externalities reflects customers' behavior when one type of product is considered superior to the other type (i.e., the vertical differentiation we consider here). As an illustration, we use the example of the Hermes silk scarf and Birkin handbag. Clearly, a Birkin handbag is considered superior to a silk scarf – one may tell it from their market prices (the price of the former is about \$10,000 while the latter is only \$400). While snobs who value uniqueness tend to show off with a Birkin handbag (and experience negative externalities when more people own Birkin(s)), their network values from the silk scarf can be considered minimal. First, the show-off value from a silk scarf is very small itself relative to that from a Birkin handbag. Second, even in the case where snobs suffer negative externalities from scarves, there are ways to reduce them largely, e.g., wrapping the silk scarf around the handle of a Birkin handbag as a personal touch instead of wearing it on the neck.¹⁸ For followers, on the contrary, mimicking other customers' purchases happen for all products.

For the market structure, we focus on the monopoly market, which we consider as a good and clean setting to focus on the interaction between externalities and bundling strategies. The competitive bundling case is beyond the scope of this paper.¹⁹

In the remaining part of this section, we first analyze the four possible selling strategies one by one and then derive the optimal strategy. For each strategy, we

¹⁸One can easily find many discussions on the social media, e.g., www.pinterest.com/pin/325244404320200046/

¹⁹The bundling decision in the competitive setting has drawn many discussions and is known as a complex problem (Armstrong and Vickers [8], Ghosh and Balachander [38], Zhou [94]). Incorporating the effects of externalities on such a case is considered to be beyond the scope of this work, and we view it as interesting future work.

directly analyze it in the cases with both positive and negative externalities – the cases without externalities are of similar structures (note that the no-externalities case can be easily obtained by setting $m = k = 0$). For the overall optimal strategy, however, we analyze the cases with and without externalities separately because, as we will see below, the structures are different with or without externalities.

We focus on the cases where externalities do not completely reverse the order of preferences, as is discussed in Section 3.4. Similar to Assumption III.3 (for products of a single category) in Section 3.4 and 3.5, we make the following assumption for each of the two product categories:

Assumption III.10. $v_S q_i - k(1 + \beta) \geq v_F q_i + m(1 + \beta), \forall i \in \{1, 2\}$.

Pure Component Strategy. Under this strategy, since products are offered separately, the monopoly firm can simply manage each product separately. For each product, the firm only needs to decide whether to sell it only to snobs or to both types of customers. For example, if the firm wants to sell product 1 only to snobs, then $D_1 = 1$, $p_1^* = V_S^1 = v_S q_1 - k$ and the revenue is $v_S q_1 - k$; if the firm wants to sell product 1 to both snobs and followers, then $D_1 = 1 + \beta$, $p_1^* = V_F^1 = v_F q_1 + m(1 + \beta)$, and the revenue is $(1 + \beta)[v_F q_1 + m(1 + \beta)]$. Thus, the revenue from selling product 1 is $R_1 = \max\{v_S q_1 - k, (1 + \beta)[v_F q_1 + m(1 + \beta)]\}$. Similarly, the revenue from selling product 2 is $R_2 = \max\{v_S q_2, (1 + \beta)[v_F q_2 + m(1 + \beta)]\}$.

Proposition III.11. The revenue of offering a pure component strategy is $\max\{v_S q_1 - k, (1 + \beta)[v_F q_1 + m(1 + \beta)]\} + \max\{v_S q_2, (1 + \beta)[v_F q_2 + m(1 + \beta)]\}$

Pure Bundling Strategy. Under this strategy, the pricing decision is straightforward – the firm decides whether to sell the bundle only to snobs, with price $p_B = v_S q_1 + v_S q_2 - k$, or to both snobs and followers, with price $p_B = \min\{v_F(q_1 + q_2) + 2m(1 + \beta)\}$.

Proposition III.12. The profit of offering a pure bundling strategy is $\max\{v_S(q_1 + q_2) - k, (1 + \beta)[v_F(q_1 + q_2) + 2m(1 + \beta)]\}$.

Total Mixed Bundling Strategy. Although we discuss the total bundling strategy here, such a form of bundling is actually very rare in the luxury product industry. The firm is effectively always offering a discount through the bundle, since customers choose the bundle only if the bundle price is smaller than the sum of the prices of both products. However, one of the key rules in the luxury retailing practice is to avoid always offering discounts (in order to maintain the brand image, Kapferer [53]).

Thanks to the representative customer model, we can naturally rule out the total mixed bundling strategy in the analysis. This is because that even if a firm offers the total mixed bundling, such a strategy effectively reduce to one of the strategies among pure component, pure bundling, and partial mixed bundling. In the total mixed bundling with positive sales of each option, only two scenarios could occur: first, snobs purchase product 1 and 2 separately, while followers purchase the bundle; second, followers purchase product 1 and 2 separately, while snobs purchase the bundle. Any other cases are irrelevant because if a customer purchases a bundle, then he does not need any product 1 or 2 anymore, and if a customer purchases a product separately, purchasing a bundle is dominated by purchasing the other products separately. Note that these two scenarios, $D_1^e = D_2^e = 1 + \beta$ always hold, as all customers effectively always purchase one unit of product 1 and product 2 no matter which option(s) they choose. However, given the fact that $p_B < p_1 + p_2$, for any type of customers, the utility is always higher from purchasing the bundle, and therefore, there is always an incentive to deviate from purchasing products individually.

We write it formally as the property below and omit the formal proof.

Proposition III.13. If a firm offers total mixed bundling, the optimal prices will never generate a scenario where there are positive sales of each option.

Partial Mixed Bundling. There are two different strategies of the partial mixed bundling: the firm can either offer both the bundle and product 1 or offer both the bundle and product 2. Recall that each customer only needs one unit of each product. Thus, if a customer purchases the bundle, then any separately-sold product 1 or 2 does not add value to him anymore. Similarly, if the customer purchases product 1 individually, then only product 2 is valuable to him in the bundle. (Obviously, this behavior is sub-optimal; instead of purchasing product 1 and the bundle, the customer can simply purchase the bundle and get the utilities from both products).

In the following analysis of the partial mixed bundling setting, we only focus on the cases where each option, i.e., the bundle and the separate product, generates positive sales. Other cases are equivalent to either the pure bundling case or the case of only selling one product.

Proposition III.14. When a firm offers partial mixed bundling, in the cases where there are positive sales of each option, the selling strategies can be of the following two forms

- If $v_F(q_1 - q_2) \leq k$, the optimal selling strategy is to offer the bundle and product 1 at prices $p_B = v_s(q_1 + q_2) - k$ and $p_1 = \min\{v_F q_1 + m(1 + \beta), v_S(q_1 + q_2) - v_F q_2 - k - m\}$. The revenue is $R = v_S(q_1 + q_2) - k(1 + \beta) + \beta \min\{v_F q_1 + m(1 + \beta), v_S(q_1 + q_2) - v_F q_2 - k - m\}$
- Otherwise, the optimal strategy is to offer the bundle and product 2 at prices $p_B = v_s(q_1 + q_2) - k$ and $p_1 = \min\{v_F q_2 + m(1 + \beta), v_S(q_1 + q_2) - v_F q_2 - k - m\}$. The revenue is $R = v_S(q_1 + q_2) - k + \beta \min\{v_F q_2 + m(1 + \beta), v_S(q_1 + q_2) - v_F q_1 - k - m\}$

For the no externalities case ($m = k = 0$), the optimal selling strategy (in the form of partial mixed bundling) can only be to offer the bundle and product 1. The intuition is that since customers have higher values for product 1, selling product 1 (together with the bundle) helps to extract more surplus.

Externalities change the structure of the partial mixed bundling strategy. When the negative externalities become large, as long as the firm still wants to use the partial mixed bundling strategy, offering product 1 separately backfires through the snob effect. In contrast, by offering product 2 separately (together with the bundle), the firm can still enjoy the follower effect.

The Optimal Strategy

We next discuss the optimal strategy considering all different types of bundling strategies. We start from the case with no externalities.

Theorem III.15. *The optimal policy without externalities are as follows:*

Under conditions $q_1 + (1 - \beta)q_2 \leq 1$ and $v_S \geq (\frac{1+\beta q_2}{q_1+q_2})v_F$, it is optimal to offer the partial mixed bundling where both the bundle and product 1 are offered; otherwise, it is optimal to offer the pure bundling.

The intuition of Theorem III.15 is straightforward – Only when the proportion of followers is large enough and snobs have similar quality sensitivity as followers, the pure bundling dominate other strategies, because it drives all followers to purchase the bundle (rather than only the product 1 separately) and such benefit from followers outweighs the loss from snobs.

We provide Figure 3.1 as a brief illustration.²⁰

Next, we discuss the case with externalities.

Theorem III.16. *The optimal policy with both positive and negative externalities is as follows:*

For the cases of $k \geq v_F(q_1 - q_2)$: it is optimal to offer the partial mixed bundling

²⁰It shows the cases of $q_1 + (1 - \beta)q_2 \leq 1$.

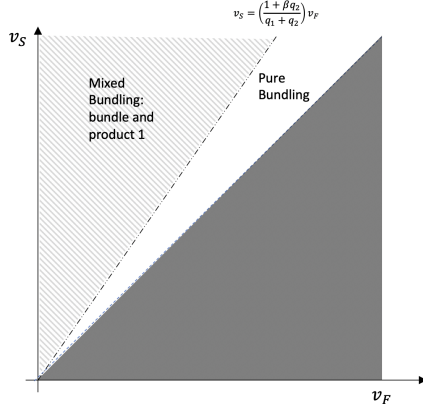


Figure 3.1: Optimal Policies without Externalities

where both the bundle and product 2 are offered if the following conditions hold.

$$(v_S - v_F)(q_1 + q_2) \geq (2 + \beta)m + k$$

$$(v_S - v_F)q_2 \geq m(1 + \beta)$$

and either $(v_S - v_F)(q_1 + q_2) - \beta v_F q_1 \geq (2 + \beta)(1 + \beta)m + k$ or $v_F q_2 + m(1 + \beta) \geq k$.

For the other cases of $k < v_F(q_1 - q_2)$: it is optimal to offer the partial mixed bundling where both the bundle and product 1 are offered if the following conditions hold.

$$(v_S - v_F)(q_1 + q_2) \geq (2 + \beta)m + k$$

$$(v_S - v_F)q_2 + \beta v_F(q_1 - q_2) \geq m(1 + \beta) + k\beta$$

and either $(v_S - v_F)(q_1 + q_2) - \beta v_F q_2 \geq (2 + \beta)(1 + \beta)m + k(1 + \beta)$ or $v_F q_1 + m(1 + \beta) \geq k$.

For all other cases, it is optimal to offer the pure bundling.

Figure 3.2 illustrates the optimal policy.²¹

²¹While the shape of the plot varies across different cases, the overall form is relatively stable. The plotted case is of conditions $q_1 + (1 - \beta)q_2 \leq 1$, $q_1 + q_2 \geq 1$, $kq_2 \leq m(1 + \beta)^2(q_1 - q_2)$, and

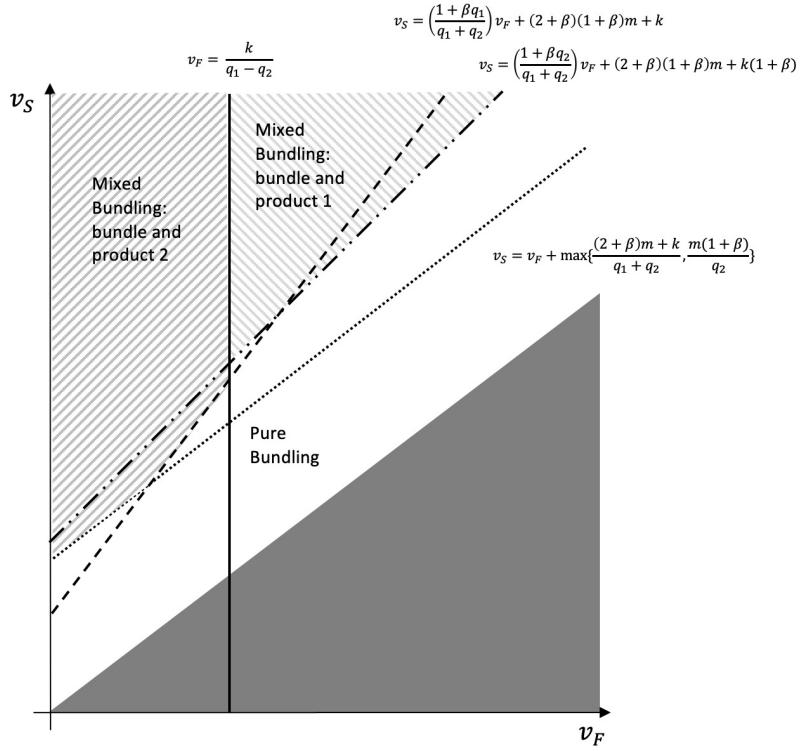


Figure 3.2: Optimal Policies with Externalities

Before jumping into the detailed explanation of Theorem III.16, it is constructive for us first to conjecture why offering the bundle and product 2 in the mixed bundling strategy would be optimal – recall that this is the counter-intuitive bundling strategy that we observe in practice. One possible conjecture would be: by only offering product 1 in the bundle, it effectively reduces the total sales of product 1 and weakens the snob effects such that the retailer can charge a higher price for the snobs. Further, by forcing any customer who purchases product 1 to also purchase product 2 (through the bundle), this strategy strengthens the follower effects such that the retailers can charge more for the followers. Although such an intuition seems plausible, we find in our analysis that only the first half of the above conjecture is accurate, but not the second half. The first half is true, especially when the sensitivity to quality of snobs (v_S) is much higher than

$$(2 + \beta)(1 + \beta)m + k \leq \frac{m(1 + \beta)^2 + k(1 + \beta)}{q_1}.$$

that of followers (v_F). The snob effect indeed determines whether to offer product 2 or product 1 separately (as is shown by the partition line $v_F = \frac{k}{q_1 - q_2}$) – when the snob effect becomes large, i.e., $v_F \leq \frac{k}{q_1 - q_2}$, the focus of the optimal policy switches from extracting more surplus from the high-value product (offering product 1 separately) to reducing the negative externalities caused by sales of the high-value product (offering product 2 separately). The second half of the above intuition, however, is wrong. It is easy to observe, from comparing Figure 3.2 with Figure 3.1, that when either the follower or the snob effect becomes stronger, the retailers are more inclined to adopt the pure bundling strategy, especially when snobs and followers have similar values for each product, i.e., the difference between v_S and v_F is small. In such cases, it is not wise for the retailer to bear losing followers on any product in order to exchange for the “exclusiveness” value for snobs (who do not have a much higher willingness-to-pay anyway). Instead, taking advantage of the follower effect on both products through pure bundling becomes more beneficial.

It is also worth noting that the pure component strategy is never the optimal strategy in any setting we consider, neither with nor without externalities. As shown in the proof of Theorem III.16, the pure component is either equivalent to the pure bundling strategy²² or strictly dominated by the mixed bundling strategy.²³ This is mainly driven by the discrete distribution of customer types.

3.7 Conclusions

Motivated by the interesting and (perhaps) counter-intuitive practices in the luxury industry, in this paper, we aim to understand the driving forces of these practices

²²It happens in two cases, one case is where the retailer only sells to snobs and the other case is where the retailer sells both products to both types of customers

²³This is the case where the retailer effectively sells both products to the snobs and only sells a single product to the followers. It can be viewed as a case of the partial mixed bundling since the realized demand is equivalent.

through the lens of externalities. Externalities of luxury products have a unique feature – They are composed of two opposite effects: snobs experience negative externalities while followers experience positive ones.

We study the joint effect of these two effects with respect to the optimal selling strategy from three aspects: (1) the product-line decision in a monopoly setting, (2) the pricing decisions in a competitive setting, and (3) the product bundling decision. In each aspect, we choose a stylized model to capture the main trade-offs and characterize the optimal decisions or the equilibrium structure. We find that snob and follower effects generally work in the same direction. In the monopoly setting, when these two externalities are not too strong, they both ease the cannibalization between products and tend to induce market segmentation (where more than one product is offered); when these two externalities become very strong, then they both tend to induce an equilibrium where only one product is offered. In the competitive setting, when it is possible for a firm to occupy the entire market, these externalities jointly induce the firm with a higher quality product to sell its product to both types of customers and to initiate a pricing war; when no firm can take the entire market, these externalities then work jointly through the two mechanisms that we identified to induce a partial segmentation in the equilibrium. For bundling decisions, the two externalities work in the same way to make the firm more inclined to choose the pure bundling (when the partial mixed bundling is still the optimal policy, the snob effect itself could change the firm's decision of which product to offer separately). In summary, we find these two opposite externalities generally work in the same direction, although the specific mechanism could be different under different market/product settings and parameters.

APPENDICES

APPENDIX A

Parameters in Simulation Experiments of Chapter I

In Section 1.5.2: $d = 5$. $T = 20$. $F_1 = ax^2 + bx + c$, where $a = 0.7$, $b = -13.5$, $c = 109.6$. $F_2 = F_1 + \beta_1$ or $F_2 = F_1 * \beta_2$, $\beta_1 \in [0, 90]$ and $\beta_2 \in [0.01, 0.9]$. α_X ($x \in \{A, B, C\}$) varies in $[0.1, 0.9]$. **In Section 1.6.2:** $d = \{4, 5, \dots, 10\}$. $T = 30$. $v(\cdot) = \gamma F(\cdot)$ where $\gamma \in (0, 1)$. $\alpha \in [0.1, 0.9]$. **In Section 1.7:** In Table 1.2 - 1.5, $F_1 = ax^2 + bx + c$, where $a = 0.7$, $b = -13.5$, $c = 109.6$, $F_2 = F_1 + \beta_1$ or $F_2 = F_1 * \beta_2$, where $\beta_1 \in [0, 90]$ and $\beta_2 \in [0.01, 0.9]$, $\alpha_X \in [0.1, 0.9]$ ($x \in \{A, B, C\}$), and $v_i(\cdot) = \gamma F_i(\cdot)$, $i \in \{1, 2\}$ where $\gamma \in [0.1, 0.9]$. $d \in \{3, 4, \dots, 10\}$ and for each d , the cost functions are discretized into d segments. Such rescaling is intended to make the performances across ds comparable.

APPENDIX B

Proof of Chapter I

Proof of Lemma I.1: (a) Suppose that we are at the beginning of period t and there are n pending orders that have not been fulfilled. If it is optimal to ship order 1, then it is also optimal to ship orders $2, 3, \dots, n$ because including orders $2, 3, \dots, n$ does not increase the current shipping cost. If, on the other hand, it is not optimal to ship order 1 in period t , then it is also not optimal to ship any subset of orders $2, 3, \dots, n$ in period t . To see this, suppose that it is optimal to ship orders $S = \{i_1, i_2, \dots, i_k\}$, where $1 < i_1 < i_2 < \dots < i_k$. Consider the following alternative shipping policy: instead of shipping S in period t , we ship them in a later period $t' < t$ when order 1 is shipped. The current shipping cost is saved and no new additional cost is incurred, which contradicts the optimality of shipping orders in S . (b) and (c) are straightforward as there are only fixed shipping costs.

Proof of Proposition I.2: For any optimal solution of $V_t(z)$, it is also feasible for $V_t(z + 1)$. Thus, $V_t(z + 1) \leq V_t(z)$. Also, extending the time-to-go horizon increases the total shipping cost, as the shipping cost is positive.

Proof of Lemma I.3: We first show that the cost-to-go function $V_t(z)$ can also be written as: $V_t(z) = \min\{F(z) + V_t(\infty), F(z - 1) + V_{t-1}(\infty), \dots, F(z - k) + V_{t-k}(\infty)\}$

where $k = \min\{t, z\} - 1$. The terms after the equality represent the costs of different alternatives. For example, $F(z) + V_t(\infty)$ is the cost of shipping all orders in period t , $F(z - 1) + V_{t-1}(\infty)$ is the cost of delaying for one period and shipping all orders in period $t - 1$, etc. It is important to note that the cost-to-go function $V_t(z)$ is completely characterized by the values of $V_t(\infty)$ for all z and t . We then introduce a technical lemma which shows that the difference between the minimum of two set of numbers is larger than the minimal pairwise difference.

LEMMA E1. *Define $x = \min\{a_1 + b_1, \dots, a_n + b_n\}$ and $y = \min\{a_1, \dots, a_n\}$. If $b_1 \geq b_2 \dots \geq b_n$, then $x - y \geq b_n$.*

Proof. Proof. Suppose that $x = a_k + b_k$ for some k . Then, $x - y \geq (a_k + b_k) - a_k = b_k \geq b_n$. \square

Then we show the proof of Lemma 2. Suppose that $t \geq d$ (the case $t < d$ can be proved in a similar manner and so is omitted). We can write: $V_t(z - 1) = \min\{F(z - 1) + V_t(\infty), F(z - 2) + V_{t-1}(\infty), \dots, F(1) + V_{t-(z-2)}(\infty)\}$ and $V_t(z) = \min\{F(z) + V_t(\infty), F(z - 1) + V_{t-1}(\infty), \dots, F(1) + V_{t-(z-1)}(\infty)\} \leq \min\{F(z) + V_t(\infty), F(z - 1) + V_{t-1}(\infty), \dots, F(2) + V_{t-(z-2)}(\infty)\} := \tilde{V}_t(z)$. Thus, $V_t(z - 1) - V_t(z) \geq V_t(z - 1) - \tilde{V}_t(z) \geq F(z - 1) - F(z)$, where the last inequality follows by Lemma E1 and the convexity of $F(\cdot)$: $a_1 = F(2) + V_{t-(z-2)}(\infty), \dots, a_n = F(z) + V_t(\infty)$, $b_1 = F(1) - F(2), \dots, b_n = F(z - 1) - F(z)$, and $F(1) - F(2) \geq F(2) - F(3) \geq \dots \geq F(z - 1) - F(z)$.

Proof of Theorem I.5: Let $X_i \sim \text{Geometric}(\alpha)$ for all i . (We assume that X_i 's are i.i.d.) Consider a sufficiently long time horizon T . If we use the same threshold τ in all periods, then the whole selling horizon can be approximately decomposed into N random cycles S_1, S_2, \dots, S_N , where $S_i = X_i + d - \tau$ and N is the smallest n such that $\sum_{i=1}^n S_i > T$. (Intuitively, $N - 1$ is the number of shipments during T periods.) Note that N is a stopping time, so by Wald's equation, $\mathbf{E}\left[\sum_{i=1}^N S_i\right] = \mathbf{E}[N] \cdot \mathbf{E}[S_1] = \mathbf{E}[N] \left(\frac{1}{\alpha} + d - \tau\right)$. Since $\sum_{i=1}^{N-1} S_i \leq T < \sum_{i=1}^N S_i$, we have $(\mathbf{E}[N] -$

1) $(\frac{1}{\alpha} + d - \tau) \leq T < \mathbf{E}[N] (\frac{1}{\alpha} + d - \tau)$. If T is sufficiently large, $\mathbf{E}[N]$ is approximately $T (\frac{1}{\alpha} + d - \tau)^{-1}$ and the average shipping costs is approximately $G_\tau(\alpha, d) := F(\tau) (\frac{1}{\alpha} + d - \tau)^{-1}$.

Proof of Proposition I.6: Denote $F(x) = a - bx$ ($F(d) \geq 0$), the long run cost is $G = \min_x G_\tau(x) = \min_x \frac{a-bx}{\frac{1}{\alpha}-1+d-x}$. We have $G'_\tau(x) = \frac{-b(\frac{1}{\alpha}-1+d-x)+(a-bx)}{(\frac{1}{\alpha}-1+d-x)^2} = \frac{b(1-\frac{1}{\alpha})+a-bd}{(\frac{1}{\alpha}-1+d-x)^2}$.
(1) When $b(1 - \frac{1}{\alpha}) + a - bd \geq 0$ ($\alpha \geq \frac{b}{a-bd+b}$), $\tau^* = 1$. (2) When $b(1 - \frac{1}{\alpha}) + a - bd < 0$ ($\alpha < \frac{b}{a-bd+b}$), $\tau^* = d - 1$.

Proof of Lemma I.7: We have $G'_\tau(\alpha) = T \frac{F'(\tau)(\frac{1}{\alpha}+d-\tau)+F(\tau)}{(\frac{1}{\alpha}+d-\tau)^2}$. As $F(z)$ is first-order continuous, G'_τ is continuous. By the property of $G'_\tau(\alpha)$, to show that τ^* decreases as α increases, we only need to show that: for any α_1 and its corresponding optimal $\tau_{\alpha_1}^*$, $G'_{\tau_{\alpha_1}^*}(\alpha_2) \geq G'_{\tau_{\alpha_1}^*}(\alpha_1)$ for any $\alpha_2 \geq \alpha_1$. As $\alpha_1 \leq \alpha_2$ and $F'(\cdot) \leq 0$, we have $F'(\tau_{\alpha_1}^*)(\frac{1}{\alpha_2} + d - \tau_{\alpha_1}^*) \geq F'(\tau_{\alpha_1}^*)(\frac{1}{\alpha_1} + d - \tau_{\alpha_1}^*)$ and $(\frac{1}{\alpha_2} + d - \tau_{\alpha_1}^*)^2 \leq (\frac{1}{\alpha_1} + d - \tau_{\alpha_1}^*)^2$. Thus, $G'_{\tau_{\alpha_1}^*}(\alpha_2) \geq G'_{\tau_{\alpha_1}^*}(\alpha_1)$. This completes the proof.

Proof of Lemma I.8: (a) For order type i ($i \in \{A, B, C\}$), suppose that we are at the beginning of period t and there are n pending orders that have not been fulfilled. If it is optimal to ship order 1, then it is also optimal to ship orders 2, 3, ..., n because including orders 2, 3, ..., n does not increase the current shipping cost. If, on the other hand, it is not optimal to ship order 1 in period t , then it is also not optimal to ship any subset of orders 2, 3, ..., n in period t . To see this, suppose that it is optimal to ship orders $S = \{i_1, i_2, \dots, i_k\}$, where $1 < i_1 < i_2 < \dots < i_k$. Consider the following alternative shipping policy: instead of shipping S in period t , we ship them in a later period $t' < t$ when order 1 is shipped. The current shipping cost is saved and no new additional cost is incurred, which contradicts the optimality of shipping orders in S .
(b) As the orders should be shipped using one shipment and the earliest order should meet the due date, the shipping cost is a function of the smallest slack time of orders.

Proof of Proposition 3 I.9: For the first part, we argue the case $z'_A \leq z_A$ in

detail. The other cases are similar. For any $V_t(z_A - 1, z_B, z_C)$, the same optimal shipping policy can be applied for $V_t(z_A, z_B, z_C)$. Thus, the optimal solution is a feasible one for $V_t(z_A, z_B, z_C)$ and $V_t(z_A - 1, z_B, z_C) \geq V_t(z_A, z_B, z_C)$. For the second part, a longer time horizon increases the total shipping cost, as the shipping cost is positive. It can be easily proved by induction. For $t = 2$, $V_1(z_A, z_B, z_C) = f(z_A, z_B, z_C)$ and $V_2(z_A, z_B, z_C) = \min_{(x_A, x_B, x_C)} f(x_A z_A, x_B z_B, x_C z_C) + \mathbf{E}[V_1(\tilde{z}_{A,\bar{x}}, \tilde{z}_{B,\bar{x}}, \tilde{z}_{C,\bar{x}})]$. Given that $f(\cdot, \cdot, \cdot) \geq 0$, $\tilde{z}_X \leq z_X$ for $X \in \{A, B, C\}$, and $V_t(z'_A, z'_B, z'_C) \geq V_t(z_A, z_B, z_C)$ for $z'_X \leq z_X$ ($X \in \{A, B, C\}$), we know that $V_1(z_A, z_B, z_C) \leq V_2(z_A, z_B, z_C)$. Suppose that it holds for all $t \leq t'$, then for $t = t' + 1$, note that $V_{t-1}(z_A, z_B, z_C) = \min_{(x_A, x_B, x_C)} f(x_A z_A, x_B z_B, x_C z_C) + \mathbf{E}[V_{t-2}(\tilde{z}_{A,\bar{x}}, \tilde{z}_{B,\bar{x}}, \tilde{z}_{C,\bar{x}})]$ and $V_t(z_A, z_B, z_C) = \min_{(x_A, x_B, x_C)} f(x_A z_A, x_B z_B, x_C z_C) + \mathbf{E}[V_{t-1}(\tilde{z}_{A,\bar{x}}, \tilde{z}_{B,\bar{x}}, \tilde{z}_{C,\bar{x}})]$. Given the induction hypothesis, it is easy to see that $V_{t-1}(z_A, z_B, z_C) \leq V_t(z_A, z_B, z_C)$. This completes the proof.

Proof of Lemma I.10: We only focus on the first part, as the second one can be argued in a similar way. Suppose that $z_A, z_B < \infty$ and it is optimal to ship product B from W1. We will argue that it is also optimal to ship order type A together with B . We divide the analysis into two cases: (1) If $z_A \leq z_B$, we prove by contradiction: Suppose the optimal policy ships only type B orders in the current period (type A order would be shipped in some later period). Consider a modified policy that does not ship type B order in the current period, but instead delays the shipment till the time when type A order would be shipped under this optimal policy. Clearly, such a modified policy would incur a lower cost. (2) If $z_A > z_B$, then shipping type A together with type B in the current period does not increase shipping cost.

Proof of Lemma I.11: We prove it by induction.

For $t = 1$: for $z_A \leq z_B$, $V_t(z_A, z_B, z_C) = \min\{F_1(\min\{z_A, z_B\}) + F_2(z_C), F_1(z_A) + F_2(\min\{z_B, z_C\})\} = F_1(z_A) + F_2(z_C)$. Thus, $V_t(z_A - 1, z_B, z_C) - V_t(z_A, z_B, z_C) = F_1(z_A -$

$1) - F_1(z_A)$. Similarly, for $z_C \leq z_B$, $V_t(z_A, z_B, z_C) = F_1(z_A) + F_2(z_C)$ and $V_t(z_A, z_B, z_C - 1) - V_t(z_A, z_B, z_C) = F_2(z_C - 1) - F_2(z_C)$; For $z_B \leq \min\{z_A, z_C\}$, $V_t(z_A, z_B, z_C) = \min\{F_1(z_B) + F_2(z_C), F_1(z_A) + F_2(z_B)\}$. Thus, $V_t(z_A, z_B - 1, z_C) - V_t(z_A, z_B, z_C) \geq \min\{F_1(z_B - 1) - F_1(z_B), F_2(z_B - 1) - F_2(z_B)\}$, where the last inequality is from Lemma E1. Then, suppose the inequalities hold for $t \leq t'$ (for some t'). Then for $t = t' + 1$, we compare the cost of all the possible shipping alternatives in $V_t(z_A, z_B, z_C)$ and $V_t(z_A, z_B, z_C)$.

We first prove the $V_t(z_A - 1, z_B, z_C) - V_t(z_A, z_B, z_C)$ part for $z_A \leq z_B$. By Lemma E1, we know that if there exist a shipping alternative, of which the corresponding costs difference between $V_t(z_A - 1, z_B, z_C)$ and $V_t(z_A, z_B, z_C)$ equals $F_1(z_A - 1) - F_1(z_A)$, while for all other shipping alternatives, the difference is larger than $F_1(z_A - 1) - F_1(z_A)$, then, then $V_t(z_A - 1, z_B, z_C) - V_t(z_A, z_B, z_C) \geq F_1(z_A - 1) - F_1(z_A)$. We classify all shipment alternatives into 2 cases: $x_A = 1$ and $x_A = \infty$. **(1)** For any given (x_A, x_B, x_C) with $x_A = 1$, when applying (x_A, x_B, x_C) , $V_t(z_A - 1, z_B, z_C) - V_t(z_A, z_B, z_C) = f(z_A - 1, z_B x_B, z_C x_C) + \mathbf{E}[V_{t-1}(\tilde{z}_A, \tilde{z}_B, \tilde{z}_C)] - f(z_A, z_B x_B, z_C x_C) - \mathbf{E}[V_{t-1}(\tilde{z}_A, \tilde{z}_B, \tilde{z}_C)] = F_1(z_A - 1) - F_1(z_A)$. The last equation is from $z_A \leq z_B$, and the fact that when it enters period $t - 1$, the slack time of type A orders, \tilde{z}_A , is the same in the cases of (z_A, z_B, z_C) and $(z_A - 1, z_B, z_C)$, as the pending type A orders are all shipped in period t (as $x_A = 1$). **(2)** For any given (x_A, x_B, x_C) with $x_A = \infty$, $V_t(z_A - 1, z_B, z_C) - V_t(z_A, z_B, z_C) = f(\infty, z_B x_B, z_C x_C) + \mathbf{E}[V_{t-1}(z_A - 2, \tilde{z}_B, \tilde{z}_C)] - f(\infty, z_B x_B, z_C x_C) - \mathbf{E}[V_{t-1}(z_A - 1, \tilde{z}_B, \tilde{z}_C)] \geq F_1(z_A - 1) - F_1(z_A)$. The inequality is from the induction hypothesis that for each scenario of \tilde{z}_B and \tilde{z}_C , $V_{t-1}(z_A - 2, \tilde{z}_B, \tilde{z}_C) - V_{t-1}(z_A - 1, \tilde{z}_B, \tilde{z}_C) \geq F_1(z_A - 2) - F_1(z_A - 1) \geq F_1(z_A - 1) - F_1(z_A)$ (the last inequality is from the convexity of F_i , $i \in \{1, 2\}$). With **(1)** and **(2)**, the $V_t(z_A - 1, z_B, z_C) - V_t(z_A, z_B, z_C)$ part naturally holds by Lemma E1.

Second, as products A and C are symmetric, $V_t(z_A, z_B, z_C - 1) - V_t(z_A, z_B, z_C) \leq$

$F_2(z_C - 1) - F_2(z_C)$ for $z_C \leq z_B$ also holds.

Then, we prove the $V_t(z_A, z_B - 1, z_C) - V_t(z_A, z_B, z_C)$ part for $z_B \leq \min\{z_A, z_C\}$. By similar logic as above, we divide the shipment alternatives into the following 2 cases.

(1) For (x_A, x_B, x_C) with $x_B = 1$, when applying (x_A, x_B, x_C) , $V_t(z_A, z_B - 1, z_C) - V_t(z_A, z_B, z_C) = f(z_A x_A, z_B - 1, z_C x_C) + \mathbf{E}[V_{t-1}(\tilde{z}_A, \tilde{z}_B, \tilde{z}_C)] - f(z_A x_A, z_B, z_C x_C) - \mathbf{E}[V_{t-1}(\tilde{z}_A, \tilde{z}_B, \tilde{z}_C)] \geq \min\{F_1(z_B - 1) - F_1(z_B), F_2(z_B - 1) - F_2(z_B)\}$. The last inequality is from $z_B \leq \min\{z_A, z_C\}$ and the fact that when it enters period $t - 1$, the slack time of type B orders, \tilde{z}_B , is the same in the cases of (z_A, z_B, z_C) and $(z_A, z_B - 1, z_C)$, as the pending type B orders are all shipped in period t (as $x_B = 1$). **(2)** For any given (x_A, x_B, x_C) with $x_B = \infty$, $V_t(z_A, z_B - 1, z_C) - V_t(z_A, z_B, z_C) = \mathbf{E}[V_{t-1}(\tilde{z}_A, z_B - 2, \tilde{z}_C)] - \mathbf{E}[V_{t-1}(\tilde{z}_A, z_B - 1, \tilde{z}_C)] \geq \min\{F_1(z_B - 1) - F_1(z_B), F_2(z_B - 1) - F_2(z_B)\}$. The inequality follows from the induction hypothesis, $V_{t-1}(\tilde{z}_A, z_B - 2, \tilde{z}_C) - V_{t-1}(\tilde{z}_A, z_B - 1, \tilde{z}_C) \geq \min\{F_1(z_B - 2) - F_1(z_B - 1), F_2(z_B - 2) - F_2(z_B - 1)\} \geq \min\{F_1(z_B - 1) - F_1(z_B), F_2(z_B - 1) - F_2(z_B)\}$ (the last inequality is from the convexity of F_i , $i \in \{1, 2\}$). With **(1)** and **(2)**, the $z_B \leq \min\{z_A, z_C\}$ part naturally holds by Lemma E1.

Proof of Theorem I.12: We first provide the following two results, Lemma E2 and Lemma E3, which will be useful in the proof of Theorem 3.

LEMMA E2. *In symmetric case where $F_1(x) = F_2(x) = F(x) \forall x$ and $\alpha_A = \alpha_C$, $V_t(z_A, \infty, \infty) - V_t(\infty, \infty, \infty) \geq F(z_A) - F(d)$, $V_t(\infty, z_B, \infty) - V_t(\infty, \infty, \infty) \geq F(z_B) - F(d)$, $V_t(\infty, \infty, z_C) - V_t(\infty, \infty, \infty) \geq F(z_C) - F(d)$. $V_t(\infty, z_B, z_C) - V_t(\infty, \infty, z_C) \geq F(z_B) - F(z_C)$, for $z_B \leq z_C$. $V_t(z_A, z_B, \infty) - V_t(z_A, \infty, \infty) \geq F(z_B) - F(z_A)$ for $z_B \leq z_A$.*

Proof of Lemma E2: We prove the claims above by induction. For $t = 1$, $V_t(z_A, \infty, \infty) - V_t(\infty, \infty, \infty) = F(z_A) \geq F(z_A) - F(d)$ and it is easy to see that other inequalities also hold. Suppose the inequalities hold for all $t \leq t_0$. Then for $t = t_0 + 1$, we show the proof for the first inequality in detail. Other inequalities can be shown in a similar

way. For the first inequality, we know that $V_t(z_A, \infty, \infty) = \min_{x_A} \{f(z_A x_A, \infty, \infty) + \mathbf{E}[V_{t-1}(\tilde{z}_{A,x_A}, \tilde{\infty}, \tilde{\infty})]\}$ and $V_t(\infty, \infty, \infty) = \mathbf{E}[V_{t-1}(\tilde{\infty}, \tilde{\infty}, \tilde{\infty})]$. To clearly show the proof, we can write down the former more explicitly: $V_t(z_A, \infty, \infty) = \min\{f(z_A, \infty, \infty) + \mathbf{E}[V_{t-1}(\tilde{\infty}, \tilde{\infty}, \tilde{\infty})], f(\infty, \infty, \infty) + \mathbf{E}[V_{t-1}(\tilde{z}_A - 1, \tilde{\infty}, \tilde{\infty})]\}$. First note that $f(z_A, \infty, \infty) = F(z_A)$ and $f(\infty, \infty, \infty) = 0$. Then, from the induction hypothesis, $\mathbf{E}[V_{t-1}(\tilde{z}_A - 1, \tilde{\infty}, \tilde{\infty})] - \mathbf{E}[V_{t-1}(\tilde{\infty}, \tilde{\infty}, \tilde{\infty})] \geq F(z_A - 1) - F(d)$. Thus, we naturally have $V_t(z_A, \infty, \infty) - V_t(\infty, \infty, \infty) = \min\{f(z_A, \infty, \infty) + \mathbf{E}[V_{t-1}(\tilde{\infty}, \tilde{\infty}, \tilde{\infty})], f(\infty, \infty, \infty) + \mathbf{E}[V_{t-1}(\tilde{z}_A - 1, \tilde{\infty}, \tilde{\infty})]\} - \mathbf{E}[V_{t-1}(\tilde{\infty}, \tilde{\infty}, \tilde{\infty})] \geq \min\{F(z_A), F(z_A - 1) - F(d)\} = F(z_A) - F(d)$, where the first inequality follows from Lemma E1.

Before stating Lemma E3, note that the DP formulation can be equivalently written as a pseudo-DP formulation as follows. Note that the shipping alternative of shipping only B from W1 or W2 is omitted by Lemma 5. For $t > 1$ and $z_A, z_B, z_C \geq 1$, we have: $V_t(z_A, z_B, z_C) = \min\{F_1(z_A) + \tilde{V}_t^1(\infty, z_B, z_C), F_1(\min\{z_A, z_B\}) + \tilde{V}_t^1(\infty, \infty, z_C), \tilde{V}_t^1(z_A, z_B, z_C)\}$, where the corresponding alternatives are “Ship A from W1,” “Ship A and B from W1,” and “Do not ship from W1,” respectively. $\tilde{V}_t^1(z_A, z_B, z_C) = \min\{F_2(z_C) + \tilde{V}_t^2(z_A, z_B, \infty), F_2(\min\{z_B, z_C\}) + \tilde{V}_t^2(z_A, \infty, \infty), \tilde{V}_t^2(z_A, z_B, z_C)\}$, where the corresponding alternatives are “Ship C from W2,” “Ship C and B from W2,” and “Do not ship from W2,” respectively. $\tilde{V}_t^2(z_A, z_B, z_C) = \mathbf{E}[V_{t-1}(g_A(z_A), g_B(z_B), g_C(z_C))]$, where $g_X(z_X)$ is a random variable which equals $z_X - 1$ with probability $1 - \alpha_X$ and d with probability α_X . For $t = 1$ and $z_A, z_B, z_C \geq 1$, $V_1(z_A, z_B, z_C) = \min\{F_1(\min\{z_A, z_B\}) + F_2(z_C), F_1(z_A) + F_2(\min\{z_B, z_C\})\}$. It is worth noting that while the notations of $\tilde{V}_t^1(z_A, z_B, z_C)$ and $\tilde{V}_t^2(z_A, z_B, z_C)$ help us to clearly define the DP formulation, it is actually mathematically equivalent to the formulation which replaces them by $V_t(z_A, z_B, z_C)$ (with a bit abuse of notation.) In this way, $V_t(z_A, z_B, z_C)$ is defined in a recursive manner. In the following proofs, we directly use this latter version for the simplicity of notation.

LEMMA E3. In a symmetric case, where $F_1(x) = F_2(x) = F(x) \forall x$ and $\alpha_A = \alpha_C$:

(1) If $z_A \geq \max\{z_B, z_C\}$, then $V_t(z_A - 1, z_B, z_C) - V_t(z_A, z_B, z_C) \geq F_1(z_A - 1) - F_1(z_A)$.

(2) If $z_C \geq \max\{z_A, z_B\}$, then $V_t(z_A, z_B, z_C - 1) - V_t(z_A, z_B, z_C) \geq F_2(z_C - 1) - F_2(z_C)$.

Proof of Lemma E3: This is an extension of Lemma 6. We prove it by induction. We show the proof of the $z_A \geq \max\{z_B, z_C\}$ part in details below. The $z_C \geq \max\{z_A, z_B\}$ part can be proved by similar logic. For $t = 1$, $V_t(z_A - 1, z_B, z_C) = F(z_A - 1) + F(z_C)$ and $V_t(z_A, z_B, z_C) = F(z_A) + F(z_C)$. The inequality naturally holds. Suppose it holds for $\forall t \leq t_0$. Then, for $t = t_0 + 1$: (1) for $z_A \geq z_B \geq z_C$, note that the shipping alternative of “ship A and B from W1” is always dominated by “ship B and C from W2.” This is because that in symmetric case $F(\min\{z_A, z_B\}) + V_t(\infty, \infty, z_C) \geq F(\min\{z_B, z_C\}) + V_t(z_A, \infty, \infty)$, as $V_t(\infty, \infty, z_C) - V_t(z_A, \infty, \infty) = V_t(z_C, \infty, \infty) - V_t(z_A, \infty, \infty) \geq F(z_C) - F(z_A) \geq F(z_C) - F(z_B)$ (where the first inequality is from Lemma E2 and the second inequality is from $z_A \geq z_B$). Second, note that “ship C from W2” cannot be optimal (it is dominated by “shipping B and C from W2”), as $F(z_C) + V_t(z_A - 1, z_B, \infty) \geq F(z_C) + V_t(z_A - 1, \infty, \infty) = F(\min\{z_B, z_C\}) + V_t(z_A - 1, \infty, \infty)$. (2) for $z_A \geq z_C \geq z_B$, “shipping C from W2” is dominated by “shipping B and C from W2” as $F(z_B) + V_t(z_A, \infty, \infty) \leq F(z_C) + V_t(z_A, z_B, \infty)$, which follows from $V_t(z_A, z_B, \infty) - V_t(z_A, \infty, \infty) \geq F(z_B) - F(z_A) \geq F(z_B) - F(z_C)$ (where the first inequality is from Lemma E2 and the second inequality is from $z_A \geq z_C$). Also, “ship A and B from W1” is dominated by “ship B and C from W2” as $F(z_B) + V_t(\infty, \infty, z_C) \geq F(z_B) + V_t(z_A, \infty, \infty)$, which follows from $V_t(\infty, \infty, z_C) = V_t(z_C, \infty, \infty) \geq V_t(z_A, \infty, \infty)$. Thus, only three shipping alternatives need to be considered in the pseudo-DP: “Ship A from W1,” “Ship C and B from W2” and “Do not ship.” $V_t(z_A, z_B, z_C) = \min\{F(z_A) + V_t(\infty, z_B, z_C), F(\min\{z_B, z_C\}) + V_t(z_A, \infty, \infty), V_{t-1}(z_A - 1, z_B - 1, z_C - 1)\}$. Then, from induction hypothesis, Lemmas E1 and E2, $V_t(z_A - 1, z_B, z_C) - V_t(z_A, z_B, z_C) \geq F(z_A - 1) - F(z_A)$. This completes the proof.

Next, we show the six-boundary result by splitting the proof into three parts. In the first part, we show the existence of thresholds $\tau_{A,t}^{AB}(z_B, z_C)$ and $\tau_{C,t}^{BC}(z_A, z_B)$. The second part is for boundaries $\tau_{A,t}^A(z_B, z_C)$ and $\tau_{C,t}^C(z_A, z_B)$, while the third part describes $\tau_{B,t}^1(z_A, z_C)$ and $\tau_{B,t}^2(z_A, z_C)$.

Part 1: The existence of thresholds $\tau_{A,t}^{AB}(z_B, z_C)$ is proved below in detail. Similar argument holds for $\tau_{C,t}^{BC}(z_A, z_B)$. We divide the proof into two cases, $z_A \leq z_B$ and $z_A > z_B$. Case 1: $z_A \leq z_B$. We show that, if for some (z_A, z_B, z_C) , the optimal policy (x_A^*, x_B^*, x_C^*) is to ship A and B from W1 ($x_A^*, x_B^* = 1$), then for $(z_A - 1, z_B, z_C)$, the optimal policy is also to ship A and B from W1. Note that the cost of applying policy (x_A^*, x_B^*, x_C^*) in state $(z_A - 1, z_B, z_C)$ is $C_{(x_A^*, x_B^*, x_C^*)}(z_A - 1, z_B, z_C) = f(z_A - 1, z_B, z_C x_C^*) + \mathbf{E}[V_{t-1}(\tilde{z}'_A, \tilde{z}_B, \tilde{z}_C)]$, where $z'_A = z_A - 1$. As $x_A^* = 1$, it is obvious that $\mathbf{E}[V_{t-1}(\tilde{z}'_A, \tilde{z}_B, \tilde{z}_C)] = \mathbf{E}[V_{t-1}(\tilde{z}_A, \tilde{z}_B, \tilde{z}_C)]$. Thus, $C_{(x_A^*, x_B^*, x_C^*)}(z_A - 1, z_B, z_C) - C_{(x_A^*, x_B^*, x_C^*)}(z_A, z_B, z_C) = f(z_A - 1, z_B, z_C x_C^*) - f(z_A, z_B, z_C x_C^*) = F(z_A - 1) - F(z_A)$. In other words, applying (x_A^*, x_B^*, x_C^*) for state $(z_A - 1, z_B, z_C)$ incurs $F(z_A - 1) - F(z_A)$ more cost than applying it for state (z_A, z_B, z_C) . Note that from Lemma 6, we know that $V_t(z_A - 1, z_B, z_C) - V_t(z_A, z_B, z_C) \geq F(z_A - 1) - F(z_A)$, which indicates that using any shipping policies for state $(z_A - 1, z_B, z_C)$ incurs at least additional cost of $F(z_A - 1) - F(z_A)$, compared to (z_A, z_B, z_C) . Thus, we can easily conclude that (x_A^*, x_B^*, x_C^*) is the optimal policy of state $(z_A - 1, z_B, z_C)$, which ships A and B from W1. Case 2: $z_A > z_B$. As (x_A^*, x_B^*, x_C^*) (ship A and B from W1) is optimal for (z_A, z_B, z_C) , then from the optimality of this policy, we know that $f(z_A, z_B, z_C x_C^*) = F(z_B) + F(z_C x_C^*)$ and $F(z_B) + F(z_C x_C^*) \leq F(z_A) + F(\min\{z_B, z_C x_C^*\})$. For $(z_A - 1, z_B, z_C)$, $C_{(x_A^*, x_B^*, x_C^*)}(z_A - 1, z_B, z_C) = f(z_A - 1, z_B, z_C x_C^*) + \mathbf{E}[V_{t-1}(\tilde{z}'_A, \tilde{z}_B, \tilde{z}_C)]$ where $f(z_A - 1, z_B, z_C x_C^*) = F(z_B) + F(z_C x_C^*)$ as $F(z_B) + F(z_C x_C^*) \leq F(z_A) + F(\min\{z_B, z_C x_C^*\}) \leq F(z_A - 1) + F(\min\{z_B, z_C x_C^*\})$. Thus, $C_{(x_A^*, x_B^*, x_C^*)}(z_A - 1, z_B, z_C) = C_{(x_A^*, x_B^*, x_C^*)}(z_A, z_B, z_C)$. Note that the shipping cost of $(z_A - 1, z_B, z_C)$ is always higher than or equal to that of (z_A, z_B, z_C) . (With $z_A > z_B$,

Lemma 6 not necessary hold. But we always have $V_t(z_A - 1, z_B, z_C) \geq V_t(z_A, z_B, z_C)$ from Proposition 3.) Thus, we can easily conclude that (x_A^*, x_B^*, x_C^*) is the optimal policy of $(z_A - 1, z_B, z_C)$, which ships A and B from W1.

Part 2: The existence of thresholds $\tau_{A,t}^1(z_B, z_C)$ is proved in detail. Similar arguments hold for $\tau_{C,t}^1(z_A, z_B)$. We want to show that if, for some $V_t(z_A, z_B, z_C)$, the optimal policy (x_A^*, x_B^*, x_C^*) is to ship both A and B from W1 or to ship only A from W1 ($x_A^* = x_B^* = 1$ or $x_A^* = 1$), then for $V_t(z_A - 1, z_B, z_C)$, the optimal policy is also to ship both A and B from W1 or to ship only A from W1. For $x_A^* = x_B^* = 1$, it is already discussed in part 1. We next prove $x_A^* = 1$ part: We want to show that if, for some $V_t(z_A, z_B, z_C)$, the optimal policy (x_A^*, x_B^*, x_C^*) is to ship only A from W1 ($x_A^* = 1$), then for $V_t(z_A - 1, z_B, z_C)$, the optimal policy is to ship both A and B from W1 or to ship only A from W1. Note that only the case of $z_A > \max\{z_B, z_C\}$ need to be analyzed in this part, because in other cases (1) and (2) listed below, shipping A from W1 incurs larger cost than shipping A and B from W1: (1) with $z_A \leq z_B$, $F(z_A) + V_t(\infty, z_B, z_C) \geq F(\min\{z_A, z_B\}) + V_t(\infty, \infty, z_C)$. (2) with $z_C \geq z_A > z_B$, $F(z_A) + V_t(\infty, z_B, z_C) \geq F(z_B) + V_t(\infty, \infty, z_C)$, as $V_t(\infty, z_B, z_C) - V_t(\infty, \infty, z_C) \geq F(z_B) - F(z_C) \geq F(z_B) - F(z_A)$ where the first inequality is from Lemma E1. Then we only need to consider the case of $z_A > \max\{z_B, z_C\}$. Again, as $x_A^* = 1$, $\mathbf{E}[V_{t-1}(\tilde{z}'_A, \tilde{z}_B, \tilde{z}_C)] = \mathbf{E}[V_{t-1}(\tilde{z}_A, \tilde{z}_B, \tilde{z}_C)]$. Thus, $C_{(x_A^*, x_B^*, x_C^*)}(z_A - 1, z_B, z_C) - C_{(x_A^*, x_B^*, x_C^*)}(z_A, z_B, z_C) = f(z_A - 1, z_B x_B^*, z_C x_C^*) - f(z_A, z_B x_B^*, z_C x_C^*) = F(z_A - 1) - F(z_A)$. From Lemma E2, we know that $V_t(z_A - 1, z_B, z_C) - V_t(z_A, z_B, z_C) \geq F(z_A - 1) - F(z_A)$, which indicates that using shipping policies other than (x_A^*, x_B^*, x_C^*) for $(z_A - 1, z_B, z_C)$ will incur cost at least $F(z_A - 1) - F(z_A)$ higher than (z_A, z_B, z_C) . Thus, we can conclude that (x_A^*, x_B^*, x_C^*) is the optimal policy of $(z_A - 1, z_B, z_C)$, which ships A from W1.

Part 3: The existence of thresholds $\tau_{B,t}^1(z_A, z_C)$ is proved in detail. Similar argu-

ments hold for $\tau_{B,t}^2(z_A, z_C)$. We want to show that if, the optimal policy (x_A^*, x_B^*, x_C^*) is to ship both A and B from W1 ($x_A^* = x_B^* = 1$), then for $V_t(z_A, z_B - 1, z_C)$, the optimal policy is also to ship both A and B from W1. Only $z_A < z_C$ need to be considered in this part. Otherwise, with $z_A \geq z_C$, shipping both A and B from W1 is dominated by shipping both B and C from W2. This is because $F(\min\{z_A, z_B\}) + V_t(\infty, \infty, z_C) \geq F(\min\{z_B, z_C\}) + V_t(z_A, \infty, \infty)$, as $V_t(\infty, \infty, z_C) - V_t(z_A, \infty, \infty) = V_t(\infty, \infty, z_C) - V_t(\infty, \infty, z_A) \geq F(z_C) - F(z_A) \geq F(\min\{z_B, z_C\}) - F(\min\{z_A, z_B\})$ (the first inequality is from Lemma E1 and the second inequality is easy to see by considering all possible relations among z_A, z_B , and z_C , satisfying the condition of $z_A \geq z_C$). We divide the proof into two cases: (1) $z_B \leq z_A < z_C$. (2) $z_A < \min\{z_B, z_C\}$. **Case 1:** $z_B \leq z_A < z_C$. $C_{(x_A^*, x_B^*, x_C^*)}(z_A, z_B - 1, z_C) - C_{(x_A^*, x_B^*, x_C^*)}(z_A, z_B, z_C) = f(z_A, z_B - 1, z_C x_C^*) - f(z_A, z_B - 1, z_C x_C^*) = F(z_B - 1) - F(z_B)$, as in symmetric case $F(\min\{z_A, z_B - 1\}) + V_t(\infty, \infty, z_C) - F(\min\{z_A, z_B\}) + V_t(\infty, \infty, z_C) = F(z_B - 1) - F(z_B)$. From Lemma 6, we know that $V_t(z_A, z_B - 1, z_C) - V_t(z_A, z_B, z_C) \geq F(z_B - 1) - F(z_B)$, which indicates that using shipping policies other than (x_A^*, x_B^*, x_C^*) for $(z_A, z_B - 1, z_C)$ incur at least additional cost of $F(z_B - 1) - F(z_B)$, compared to (z_A, z_B, z_C) . Thus, we can conclude that (x_A^*, x_B^*, x_C^*) is the optimal policy of $(z_A, z_B - 1, z_C)$, which ships A and B from W1. **Case 2:** $z_A < \min\{z_B, z_C\}$. $C_{(x_A^*, x_B^*, x_C^*)}(z_A, z_B - 1, z_C) - C_{(x_A^*, x_B^*, x_C^*)}(z_A, z_B, z_C) = f(z_A, z_B - 1, z_C x_C^*) - f(z_A, z_B - 1, z_C x_C^*) = 0$. From Proposition 3, we know that $V_t(z_A, z_B - 1, z_C) \geq V_t(z_A, z_B, z_C)$, which indicates that for $(z_A, z_B - 1, z_C)$ using shipping policies other than (x_A^*, x_B^*, x_C^*) incur higher cost than (z_A, z_B, z_C) . Thus, (x_A^*, x_B^*, x_C^*) is the optimal policy of $(z_A, z_B - 1, z_C)$, which ships A and B from W1.

Proof of Lemma I.13: First, we show a result of the effect of α on $V_t(z_A, z_B, z_C) - V_{t-1}(z_A, z_B, z_C)$, which will be useful in the proof of Lemma 7. Denote the cost-to-go function as $\bar{V}_t(z_A, z_B, z_C)$ and $V_t(z_A, z_B, z_C)$ for $\bar{\alpha}$ and α , respectively. Note that under optimal stationary policy, $V_t(z_A, z_B, z_C) - V_{t-1}(z_A, z_B, z_C)$ converges to the

expected one-period cost. For $\bar{\alpha} \geq \alpha$, the expected one-period cost of $\bar{\alpha}$ is larger than that of α . This is because on average, more orders arrive per period in the case of $\bar{\alpha}$, which incurs more shipping cost. Thus, $\bar{V}_t(z_A, z_B, z_C) - \bar{V}_{t-1}(z_A, z_B, z_C) \geq V_t(z_A, z_B, z_C) - V_{t-1}(z_A, z_B, z_C)$. Next, we show the proof of $\tau_A^A(z_B, z_C)$ in detail. Other cases are similar. For any z_B, z_C and $z_A = \tau_A^A\{z_B, z_C\}$, $V_t(z_A, z_B, z_C) = F_1(z_A) + V_t(\infty, z_B, z_C) \leq V_{t-1}(z_A - 1, z_B - 1, z_C - 1)$. Then consider $\bar{V}_t(z_A, z_B, z_C)$. Following $\bar{V}_t(z_A, z_B, z_C) - \bar{V}_{t-1}(z_A, z_B, z_C) \geq V_t(z_A, z_B, z_C) - V_{t-1}(z_A, z_B, z_C)$ and $\bar{V}_1(z_A, z_B, z_C) = V_1(z_A, z_B, z_C)$ ($\forall z_A, z_B, z_C$), we have $\bar{V}_t(\infty, z_B, z_C) - V_t(\infty, z_B, z_C) \geq \bar{V}_{t-1}(z_A - 1, z_B - 1, z_C - 1) - V_{t-1}(z_A - 1, z_B - 1, z_C - 1)$. Thus, shipping only A (the policy corresponding to $\tau_A^A\{z_B, z_C\}$) is not necessarily the policy with smallest cost, which indicates the decreasing of $\tau_A(z_B, z_C)$.

Proof of Lemma I.14: For the first part: suppose order z_i is shipped and n orders are shipped in the current period. Suppose it is optimal to ship z_k ($k > i + n - 1$). Then there must be some order z_j ($i \leq j \leq i + n - 1$) shipped in the future. By exchanging z_k and z_j , the shipping cost in the current period can be kept the same while the future cost is reduced. This contradicts with the optimality.

For the second part: suppose a package (denote P_i) includes z_i as the most urgent order and n orders are shipped in it. Suppose it is optimal to include z_k ($k > i + n - 1$) in the current package. Then there must be some order z_j ($i \leq j \leq i + n - 1$) shipped in another package (denote P_j). By exchanging z_k and z_j , the shipping cost of package P_j is reduced while the cost of package P_i does not change. This contradicts with the optimality.

Proof of Theorem I.15: Before we show the proof, we first introduce notations for the changes of states across periods. For any state i in period t , we denote its predecessor in period $t + 1$ as $p(i)$ and successors in period $t - 1$ as $s(i)$. We start from proving parts 1 and 2. The states for $d = 3$ are (∞) , (1) , (2) , (3) , $(1, 2)$, $(1, 3)$, $(2, 3)$

and $(1, 2, 3)$. We prove, by induction, that in $(1, 2)$, $(1, 3)$, $(2, 3)$ and $(1, 2, 3)$, it is not optimal to ship only partial of the orders. For $t = 1$, by definition, it is optimal to ship all of the pending orders. Suppose it is not optimal to ship only partial of the orders for $t \leq t' - 1$ periods. In other words, it is optimal to either ship all orders or do not ship in states $(1, 3)$, $(1, 2)$, $(1, 2, 3)$ and $(2, 3)$. We only list these four cases, as in other cases it is natural to either ship the order or not. For $t = t'$, we consider the following 4 cases. **Case 1: State $(1, 3)$.** Suppose in period t under the optimal policy π^* , it is optimal to ship only order 1 in state $(1, 3)$. Note that $p(1, 3) = (2)$ in period $t + 1$ and the optimal policy for (2) has to be “Do not ship.” We next argue that “Do not Ship” cannot be an optimal policy for (2) . If (2) is not shipped in period $t + 1$, the state in period t becomes (1) with probability α and $(1, 3)$ with probability $1 - \alpha$. In either case, order 1 is shipped, resulting in a higher cost than in the previous period. Also, if the optimal policy of $(1, 3)$ is to ship both orders but in separate packages, the above argument also holds. Thus, if state $(1, 3)$ is reachable, the optimal policy has to be shipping both orders in one package: $V_t(1, 3) = C(1, 3) + V_t(\infty) \leq C(1) + V_t(3)$, where $C(1, 3) = F(1) + 2v(1) \leq F(1) + v(1) + F(3) + v(3)$, and the following inequality holds, $v(1) \leq F(3) + v(3)$. Note that the equation $(v(1) \leq F(3) + v(3))$ does not impose any assumption about the relation between $F(\cdot)$ and $v(\cdot)$. It only states that, if state $(1, 3)$ is reachable, then this equation must hold. In other words, if this equation does not hold, then state $(1, 3)$ cannot be reached under the optimal policy. **Case 2: State $(1, 2)$.** Suppose that under the optimal policy π^* , it is optimal to ship only order 1 in state $(1, 2)$ in period t . Then, in period $t - 1$, the remaining order (2) has successor $s(2) = \{(1, 3), (1)\}$ with probability α and $1 - \alpha$, respectively. Thus, state $(1, 3)$ is reachable and inequality of $v(1) \leq F(3) + v(3)$ holds. Then, it can be shown that both orders in state $(1, 2)$ need to be shipped in one package: $C(1, 2) = \min\{F(1) + 2v(1), F(1) + v(1) + F(2) + v(2)\} = F(1) + 2v(1) = C(1, 3)$, as $F(1) + 2v(1) \leq$

$F(1) + v(1) + F(3) + v(3) \leq F(1) + v(1) + F(2) + v(2)$, where the first inequality follows from $v(1) \leq F(3) + v(3)$. Also, note that $C(1, 2) + V_t(\infty) \leq C(1) + V_t(2)$, as $C(1, 2) + V_t(\infty) = C(1, 3) + V_t(\infty) \leq C(1) + V_t(3) \leq C(1) + V_t(2)$. Thus, shipping both order 1 and 2 in one package incurs smaller cost than π^* , which contradicts the optimality of π^* . Thus, if state $(1, 2)$ is reachable, the optimal policy is to ship both orders in one package. $V_t(1, 2) = C(1, 2) + V_t(\infty)$, where $C(1, 2) = F(1) + 2v(1)$. **Case 3: State $(1, 2, 3)$.** There are 4 shipment alternatives for state $(1, 2, 3)$: ship 1 alone, ship 1 and 2, ship 1 and 3 and ship all orders. We argue that the first three alternatives cannot be optimal. First, suppose that under the optimal policy π^* it is optimal to ship only order $(1, 2)$ in state $(1, 2, 3)$ in period t . Note $p(1, 2, 3) = (2, 3)$ where the optimal policy for $(2, 3)$ should be "Do not ship." And $s(2, 3) = \{(1, 2, 3), (1, 2)\}$ with probability α and $1 - \alpha$, respectively. As $(1, 2)$ is reachable, from the results in Case 2, $C(1, 2) = F(1) + 2v(1)$. Note that in both cases of $(1, 2, 3)$ and $(1, 2)$, orders $(1, 2)$ are shipped in one package and incur cost $F(1) + 2v(1)$. Consider policy $\tilde{\pi}$ which chooses to ship in state $(2, 3)$ in period $t + 1$ and keeps other decisions the same as π^* . The cost of $\tilde{\pi}$ is $F(2) + 2v(2)$, which is smaller than that of π^* , which contradicts the optimality of π^* . Also, note that, if the optimal policy for $(1, 2, 3)$ is to ship order 1 and 2 in package one and order 3 in package two, the argument also holds. Second, suppose it is optimal to ship 1 and 3 in state $(1, 2, 3)$ in period t . If (2) is not shipped in period t , the state in period $t - 1$ becomes $(1, 3)$ with probability α and (1) with probability $1 - \alpha$. As $(1, 3)$ is reachable, $v(1) \leq F(3) + v(3)$ holds. Then, $C(1, 3) = F(1) + 2v(1) \leq F(1) + v(1) + F(3) + v(3)$. Thus, 1 and 3 are shipped in the same package. Consider policy $\tilde{\pi}$ that ships orders 1 and 2 in period t and ships order 3 following the same policy as order 2 in policy π^* . $\tilde{\pi}$ incurs the same cost in period t but lower cost in future period, which contradicts the optimality of policy π^* . Third, suppose under the optimal policy π^* it is optimal to ship only

order (1) in state (1, 2, 3) in period t . Thus, in period $t - 1$, the remaining orders (2, 3) becomes $s(2, 3) = \{(1, 2, 3), (1, 2)\}$ with probability α and $1 - \alpha$, respectively. From induction hypothesis, it is optimal to ship all orders in state (1, 2, 3) and (1, 2) in period $t - 1$. Thus, the remaining order 2 in period t incurs at least variable cost $v(1)$ in period $t - 1$. Consider policy $\tilde{\pi}$ which chooses to ship (1, 2) in state (1, 2, 3) in period t , and keeps other decisions the same as π^* . $\tilde{\pi}$ incurs additional cost of $v(1)$ in period t , while the cost decreases at least by $v(1)$ in period $t - 1$. It contradicts the optimality of policy π^* . Thus, if state (1, 2, 3) is reachable, the optimal policy is to ship all orders. Whether to ship them in one package or in separate packages depends on the relationship of $F(\cdot)$ and $v(\cdot)$: $C(1, 2, 3) = \min\{F(1) + 3v(1), F(1) + v(1) + F(2) + 2v(2), F(1) + v(1) + F(2) + v(2) + F(3) + v(3)\}$. **Case 4: State (2, 3).** We show that shipping either (2) or (3) along can not be optimal. First, suppose under the optimal policy π^* it is optimal to ship only order 2 in state (2, 3) in period t . In period $t - 1$, the remaining order (3) has successor $s(3) = \{(2, 3), (2)\}$ with probability α and $(1 - \alpha)$, respectively. From induction hypothesis, in state (2, 3), it is optimal to either ship both orders or not to ship. Thus, the remaining order (3) in period t will be shipped with other orders arriving in later periods, which incurs at least variable cost $v(2)$. Consider policy $\tilde{\pi}$ which ships (2, 3) in period t and keep other decisions the same as in policy π^* . $\tilde{\pi}$ incurs $v(2)$ higher cost in period t while decrease at least $v(2)$ cost in later periods. It contradicts the optimality of policy π^* . Second, suppose under the optimal policy π^* it is optimal to ship only order 3 in state (2, 3) in period t . In period $t - 1$, the remaining order (2) has successor $s(2) = \{(1, 3), (1)\}$ with probability α and $(1 - \alpha)$, respectively. It is easy to see that exchange the policy of order 2 and 3 in period t incurs $v(2) - v(3)$ higher cost in period t while decreases at least $v(2) - v(1)$ cost in period $t - 1$. It contradicts with the optimality of π^* . Thus, if state (2, 3) is reachable, the optimal policy is to ship both orders. Whether to ship

them in one package or in separate packages depends on the relation between $F(\cdot)$ and $v(\cdot)$: $C(2, 3) = \min\{F(2) + 2v(2), F(2) + v(2) + F(3) + v(3)\}$.

Finally, we prove the third part of Theorem 4 using the result from the first and second parts of Theorem 4. Note that it is sufficient to consider the following three scenarios: Scenario 1, if it is optimal to ship for state (1), then it is optimal to ship in states (1, 2) and (1, 3); Scenario 2, if it is optimal to ship for state (1, 2) or (1, 3), then it is optimal to ship in state (1, 2, 3); Scenario 3, if it is optimal to ship for state (2), then it is optimal to ship in state (2, 3). The first two scenarios are obvious, as the order needs to be shipped for $z_1 = 1$. Thus, we only need to prove the third scenario. As the optimal policy for $V_t(2)$ ($\forall t > 1$) is to ship, $F(2) + v(2) + V_t(\infty) \leq \alpha V_{t-1}(1, 3) + (1 - \alpha)V_{t-1}(1)$. Then, for state (2, 3), $F(2) + 2v(2) + V_t(\infty) \leq \alpha V_{t-1}(1, 3) + (1 - \alpha)V_{t-1}(1) + v(2) \leq \alpha V_{t-1}(1, 2, 3) + (1 - \alpha)V_{t-1}(1, 2)$. Thus, it is also optimal to ship (all pending orders) in state (2, 3).

Proof of Theorem I.16: We first want to show the cases where apply varying threshold $\tau(m)$ and constant threshold τ are equivalent. First, note that for $\tau(m)$ where $m = 3$, there is only one case where three orders are accumulated, i.e. $(z_1, z_2, z_3) = (1, 2, 3)$. As the orders must be shipped when the slack time equals 1, setting any value to $\tau(3)$ has the same impact on the policy. Thus, we only need to focus on the values of $\tau(1)$ and $\tau(2)$. Second, note that for $\tau(m)$ where $m = 2$, setting $\tau(2) = 3$ is equivalent to $\tau(2) = 2$, as the smallest slack time of the two orders cannot be larger than 2. Thus, the only varying thresholds, which are not equivalent to a constant threshold, are $\tau(1) = 1$ and $\tau(2) = 2$. We denote the average cost per period incurred by these threshold as $C_{1,2}$, and the average cost incurred by constant threshold $\tau = m$ by C_m .

Thus, the gap between varying threshold and constant threshold is equivalent to the gap between $C_{1,2}$ and $\min_{m \in \{1,2,3\}} C_m$. By simple Markov Chain argument, it is easy to see that for policies where orders are shipped in one package, $C_1 =$

$$F(1)\frac{\alpha}{1+3\alpha} + v(1)\frac{\alpha(1+2\alpha)}{1+3\alpha}, C_2 = F(2)\frac{\alpha}{1+2\alpha} + v(2)\frac{\alpha(1+\alpha)}{1+2\alpha}, \text{ and } C_{1,2} = [F(2)+2v(2)]\frac{\alpha^2}{1+3\alpha-\alpha^2} + F(1)\frac{\alpha(1-\alpha)}{1+3\alpha-\alpha^2} + v(1)\frac{\alpha(1-\alpha^2)}{1+3\alpha-\alpha^2}.$$

Then, for $\beta = \frac{(1+3\alpha)(1-\alpha)}{1+3\alpha-\alpha^2}$. $\min_{m \in \{1,2,3\}} C_m \leq \beta C_1 + (1-\beta)C_2 = C_{1,2} + [v(1) - v(2)]\frac{\alpha^2(1-\alpha)}{1+3\alpha-\alpha^2}$. Thus, $\frac{\min_{m \in \{1,2,3\}} C_m - C_{1,2}}{C_{1,2}} \leq \frac{[v(1)-v(2)]\frac{\alpha^2(1-\alpha)}{1+3\alpha-\alpha^2}}{[F(2)+2v(2)]\frac{\alpha^2}{1+3\alpha-\alpha^2} + F(1)\frac{\alpha(1-\alpha)}{1+3\alpha-\alpha^2} + v(1)\frac{\alpha(1-\alpha^2)}{1+3\alpha-\alpha^2}} = \frac{[v(1)-v(2)]\frac{\alpha^2(1-\alpha)}{1+3\alpha-\alpha^2}}{[F(2)+2v(2)]\frac{1}{1-\alpha} + F(1)\frac{1}{\alpha} + v(1)(1+\frac{1}{\alpha})} \leq \frac{\Delta}{a_1\gamma\Delta + (a_1+a_2)v(2)}$, where $a_1 = \frac{2}{\alpha} + 1$ and $a_2 = \frac{3}{1-\alpha}$. The bounds above goes to 0 when $\Delta \rightarrow 0$, $\alpha \rightarrow 1$, or $\alpha \rightarrow 0$.

For the other extreme case where Δ goes to ∞ . It is obvious that C_1 cannot be optimal. Thus, $\frac{\min_{m \in \{1,2,3\}} C_m - C_{1,2}}{C_{1,2}} \leq \frac{C_2 - C_{1,2}}{C_{1,2}}$. As $C_2 - C_{1,2} = F(2)\frac{\alpha}{1+2\alpha} + v(2)\frac{\alpha(1+\alpha)}{1+2\alpha} - [F(2)+2v(2)]\frac{\alpha^2}{1+3\alpha-\alpha^2} - F(1)\frac{\alpha(1-\alpha)}{1+3\alpha-\alpha^2} - v(1)\frac{\alpha(1-\alpha^2)}{1+3\alpha-\alpha^2} = F(2)[\frac{\alpha}{1+2\alpha} - \frac{\alpha^2}{1+3\alpha-\alpha^2}] + v(2)[\frac{\alpha(1+\alpha)}{1+2\alpha} - \frac{2\alpha^2}{1+3\alpha-\alpha^2}] - F(1)\frac{\alpha(1-\alpha)}{1+3\alpha-\alpha^2} - v(1)\frac{\alpha(1-\alpha^2)}{1+3\alpha-\alpha^2} \leq \frac{1}{1+3\alpha-\alpha^2} [\frac{\alpha^2(1-\alpha)^2}{1+2\alpha} F(1) - \frac{\alpha(1-\alpha)(1+3\alpha)}{1+2\alpha} \Delta]$, we have $\frac{C_2 - C_{1,2}}{C_{1,2}} \leq \frac{F(2)\frac{1+(1+3\alpha)(3-\alpha)}{1+2\alpha} - \Delta(\frac{1+\gamma}{\alpha} + \gamma)}{[F(2)+2v(2)]\frac{1}{1-\alpha} + F(1)\frac{1}{\alpha} + v(1)(1+\frac{1}{\alpha})} \leq \frac{F(2)\frac{1+(1+3\alpha)(3-\alpha)}{1+2\alpha} - \Delta(\frac{1+\gamma}{\alpha} + \gamma)}{F(2)\frac{1}{\alpha(1-\alpha)} + \Delta\frac{1}{\alpha}}$. It goes to negative with $\Delta \rightarrow \infty$. Note that this is quite intuitive: as $\Delta \rightarrow \infty$, it is not optimal to set any threshold to 1. Thus, it is easy to see that C_2 should has smaller cost than $C_{1,2}$. And note that in such cases, the optimal policy with varying threshold (including scenarios of constant thresholds) should be the same as the optimal policy with constant threshold. In other words, the actual bound in this case becomes $\frac{\tilde{C} - C_0}{C_0} \leq \max\{\frac{F(2)\frac{1+(1+3\alpha)(3-\alpha)}{1+2\alpha} - \Delta(\frac{1+\gamma}{\alpha} + \gamma)}{F(2)\frac{1}{\alpha(1-\alpha)} + \Delta\frac{1}{\alpha}}, 0\}$. We can further loose this bound and write it compactly: $\frac{\tilde{C} - C_0}{C_0} \leq \frac{F(2)\frac{1+(1+3\alpha)(3-\alpha)}{1+2\alpha}}{F(2)\frac{1}{\alpha(1-\alpha)} + \Delta\frac{1}{\alpha}}$, which goes to 0 as $\Delta \rightarrow \infty$.

Proof of Theorem I.17: Denote the packages shipped under any policy π in T periods as p_i^π ($i \leq k^\pi$, where k^π is the total number of packages). In package p_i^π , denote the smallest slack time as $z(p_i^\pi)$ and the number of orders in the package as $m(p_i^\pi)$. The total cost of any policy π is $E[\sum_{i \leq k^\pi} F(z(p_i^\pi)) + m(p_i^\pi) v(z(p_i^\pi))]$. **For the first bound:** We derive the relation between the costs C^* of the optimal policy π^* and the costs C_F of the policy π_F which considers only fixed cost. $C^* = E[\sum_{i \leq k^{\pi^*}} F(z(p_i^{\pi^*})) + m(p_i^{\pi^*}) v(z(p_i^{\pi^*}))] \geq E[\sum_{i \leq k^{\pi^*}} F(z(p_i^{\pi^*})) + \gamma F(z(p_i^{\pi^*}))] = E[(1 + \gamma) \sum_{i \leq k^{\pi^*}} F(z(p_i^{\pi^*}))] \geq (1 + \gamma)E[\sum_{i \leq k^{\pi_F}} F(z(p_i^{\pi_F}))]$. The first inequality

follows from $v = \gamma F$ and the fact that there must be at least one order in each packages. The second inequality follows from the optimality of π_F , which considers only fixed cost. Thus, $C^* \geq (1 + \gamma)C_F$. Let $C(\pi_F)$ denote the cost of applying π_F in the case with both fixed and variable cost. Then, $C(\pi_F) \geq C^* \geq (1 + \gamma)C_F$. As π_F is one-threshold(τ) policy, the cost of $C(\pi_F)$ and C_F can be derived explicitly as $C_F = \frac{T}{\frac{1}{\alpha} + d - \tau} F(\tau)$ and $C(\pi_F) = \frac{T}{\frac{1}{\alpha} + d - \tau} [F(\tau) + \mathbb{E}(m)v(\tau)]$ where $\mathbb{E}(m) = 1 + (d - \tau)\alpha$. Plugging in the inequality $C(\pi_F) \geq C^* \geq (1 + \gamma)C_F$ and by simple algebra, $\frac{C(\pi_F) - C^*}{C^*} \leq \frac{1}{1 + \gamma} \frac{C(\pi_F)}{C_F} - 1 \leq \frac{\gamma(d - \tau)\alpha}{1 + \gamma}$. **For the second bound:** We derive the relation between C^* and the costs C_v of the policy π_v which considers only variable cost. $C^* = E[\sum_{i \leq k\pi^*} F(z(p_i^{\pi^*})) + m(p_i^{\pi^*})v(z(p_i^{\pi^*}))] = E[\sum_{i \leq k\pi^*} \frac{1}{\gamma} v(z(p_i^{\pi^*})) + m(p_i^{\pi^*})v(z(p_i^{\pi^*}))] \geq E[\sum_{i \leq k\pi^*} \frac{1}{\gamma} v(d) + m(p_i^{\pi^*})v(d)] \geq \frac{v(d)}{\gamma} \frac{T}{\frac{1}{\alpha} + d} + C_v$. The last inequality follows from the fact that the expected number of packages shipped in the long run is $\frac{T}{\frac{1}{\alpha} + d}$ (when always arrange shipment corresponding to the shipping method of cost $v(d)$) and the optimality of C_v , which considers only variable cost. Denote the cost of applying π_v to the case with both fixed and variable cost as $C(\pi_v)$. Then, $C(\pi_v) \geq C^* \geq C_v + \frac{v(d)}{\gamma} \frac{T}{\frac{1}{\alpha} + d}$. Note that $C(\pi_v)$ and C_v can be written explicitly as orders are shipped upon arrivals: $C(\pi_v) = \frac{T}{\alpha} [F(d) + v(d)] = T\alpha[F(d) + v(d)]$ and $C_v = T\alpha v(d)$. Plugging $C(\pi_v)$ and C_v into the above inequality and by simple algebra, $\frac{C(\pi_v) - C^*}{C^*} \leq \frac{C(\pi_v) - C_v - \frac{v(d)}{\gamma} \frac{T}{\frac{1}{\alpha} + d}}{C_v + \frac{v(d)}{\gamma} \frac{T}{\frac{1}{\alpha} + d}} = \frac{d\alpha}{\gamma(1 + d\alpha) + 1}$.

APPENDIX C

Proof of Chapter II

Proof of Lemma II.1. Note that (2.5) can be written more compactly as follows:

$$\begin{aligned} D^{*,b} = \quad & \min_x d'x \\ & \text{s.t. } Ax \geq b, x \geq 0 \end{aligned}$$

where $d = (d_t)$ with $d_t = c + h \cdot (T - t + 1)$ (we simply replace the variable y_t in the original formulation with $\sum_{s=1}^{t-1} x_{s-L} - \sum_{s=1}^{t-1} \mu_s$), $b = (b_t)$ with $b_t = \sum_{s=1}^t \mu_s + \beta_1^t(\alpha)$ for $1 \leq t \leq L + 1$ and $b_t = \sum_{s=1}^t \mu_s + \beta_{t-L}^{L+1}(\alpha)$ for $L + 2 \leq t \leq T$, and A is a $T \times T$ lower-triangular matrix defined below:

$$A = \begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ 1 & 1 & 0 & \cdots & 0 \\ 1 & 1 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & 1 & \cdots & 1 \end{bmatrix}.$$

The recursive formulas stated in Lemma II.1 immediately follow from the fact that d_t is decreasing in t and A is a lower-triangular matrix whose lower-triangle components are all equal to one (i.e., given $x_{1-L}, \dots, x_{t-1-L}$, we always want to set x_{t-L} to be as small as possible). ■

Proof of Lemma II.2. Let $I_t^{*,b}$ and $x_t^{*,b}$ denote the inventory level and ordering decision at period t under an optimal control in the backorder system. First note that, by Jensen's inequality,

$$\begin{aligned}
C^{*,b} &= \sum_{t=1}^T c \cdot \mathbf{E}[x_{t-L}^{*,b}] + \sum_{t=1}^T h \cdot \mathbf{E}[(I_t^{*,b} + x_{t-L}^{*,b} - D_t)^+] + \sum_{t=1}^T p \cdot \mathbf{E}[(D_t - I_t^{*,b} - x_{t-L}^{*,b})^+] \\
&\geq \sum_{t=1}^T c \cdot \mathbf{E}[x_{t-L}^{*,b}] + \sum_{t=1}^T h \cdot \left(\mathbf{E}[I_t^{*,b}] + \mathbf{E}[x_{t-L}^{*,b}] - \mu_t \right)^+ + \sum_{t=1}^T p \cdot \left(\mu_t - \mathbf{E}[I_t^{*,b}] - \mathbf{E}[x_{t-L}^{*,b}] \right)^+ \\
&:= \Phi.
\end{aligned}$$

By (2.2), (2.3) and (2.4), the terms $\mathbf{E}[I_t^{*,b}]$ and $\mathbf{E}[x_t^{*,b}]$ satisfy:

$$\begin{aligned}
\sum_{s=1}^t \mathbf{E}[x_{s-L}^{*,b}] - \sum_{s=1}^t \mu_s &\geq \beta_1^t(\alpha), \quad \forall 1 \leq t \leq L+1, \\
\sum_{s=1}^t \mathbf{E}[x_{s-L}^{*,b}] - \sum_{s=t-L}^t \mu_s &\geq \beta_{t-L}^{L+1}(\alpha), \quad \forall L+2 \leq t \leq T.
\end{aligned}$$

Therefore, the following is a feasible solution to **DET**:

$$\begin{aligned}
y_t &= \mathbf{E}[I_t^{*,b}], \\
x_t &= \mathbf{E}[x_t^{*,b}], \\
z_{t+1} &= (\mathbf{E}[I_t^{*,b}] + \mathbf{E}[x_{t-L}^{*,b}] - \mu_t)^+ = (y_{t+1})^+, \\
m_{t+1} &= (\mu_t - \mathbf{E}[I_t^{*,b}] - \mathbf{E}[x_{t-L}^{*,b}])^+ = (-y_{t+1})^+.
\end{aligned}$$

We conclude that $C^{*,b} \geq \Phi \geq D^{*,b}$. ■

Proof of Lemma II.4. For notational brevity, we will suppress the notational dependency on α and simply write $x_t^{H_b,b}(\alpha)$ and $x_t^{D,b}(\alpha)$ as $x_t^{H_b,b}$ and $x_t^{D,b}$, respectively. We will prove the lemma by induction.

For $t \leq 1$, by definition of H_b , $x_t^{H_b,b} = x_t^{D,b}$ and, therefore, $x_t^{H_b,b} - x_t^{D,b} = 0$.

For $t = 2$, by definition of H_b , we have:

$$\begin{aligned} x_2^{H_b,b} &= \left(y_2^{D,b} + \sum_{s=2-L}^2 x_s^{D,b} - I_2^{H_b,b} - \sum_{s=2-L}^1 x_s^{H_b,b} \right)^+ \\ &= (x_2^{D,b} + D_1 - \mu_1)^+ \\ &= (x_2^{D,b} + D_1 - \mu_1) + (\mu_1 - x_2^{D,b} - D_1)^+. \end{aligned}$$

So, we can bound:

$$\begin{aligned} \mathbf{E} \left[x_2^{H_b,b} - x_2^{D,b} \right] &= \mathbf{E} \left[(D_1 - \mu_1) + (\mu_1 - x_2^{D,b} - D_1)^+ \right] \\ &\leq (\mu_1 - x_2^{D,b} - D_1)^+. \end{aligned}$$

Now, for any $t \geq 2$, by definition of H_b , we have:

$$\begin{aligned} x_t^{H_b,b} &= \left(\sum_{s=2}^t x_s^{D,b} - \sum_{s=1}^{t-1} \mu_s - \sum_{s=2}^{t-1} x_s^{H_b,b} + \sum_{s=1}^{t-1} D_s \right)^+ \\ &= \sum_{s=2}^t x_s^{D,b} - \sum_{s=1}^{t-1} \mu_s - \sum_{s=2}^{t-1} x_s^{H_b,b} + \sum_{s=1}^{t-1} D_s + \left(\sum_{s=1}^{t-1} \mu_s + \sum_{s=2}^{t-1} x_s^{H_b,b} - \sum_{s=1}^{t-1} D_s - \sum_{s=2}^t x_s^{D,b} \right)^+ \end{aligned}$$

and, therefore,

$$\sum_{s=2}^t (x_s^{H_b,b} - x_s^{D,b}) = \sum_{s=1}^{t-1} (D_s - \mu_s) + \left(\sum_{s=1}^{t-1} \mu_s + \sum_{s=2}^{t-1} x_s^{H_b,b} - \sum_{s=1}^{t-1} D_s - \sum_{s=2}^t x_s^{D,b} \right)^+ \quad (\text{C.1})$$

Suppose that the following inequality holds for $t \leq t'$:

$$\mathbf{E} \left[\sum_{s=2}^t (x_s^{H_b,b} - x_s^{D,b}) \right] \leq \mathbf{E} \left[\sum_{s=1}^{t-1} (\mu_s - x_{s+1}^{D,b} - D_s)^+ \right]. \quad (\text{C.2})$$

Then, for $t = t' + 1$, we have:

$$\begin{aligned} \mathbf{E} \left[\sum_{s=2}^t (x_s^{H_b,b} - x_s^{D,b}) \right] &= \mathbf{E} \left[\sum_{s=1}^{t-1} (D_s - \mu_s) \right. \\ &\quad \left. + \left(\sum_{s=1}^{t-1} \mu_s + \sum_{s=2}^{t-1} x_s^{H_b,b} - \sum_{s=1}^{t-1} D_s - \sum_{s=2}^t x_s^{D,b} \right)^+ \right] \\ &= \mathbf{E} \left[\left(\sum_{s=1}^{t-2} \mu_s + \sum_{s=2}^{t-1} x_s^{H_b,b} - \sum_{s=1}^{t-2} D_s - \sum_{s=2}^{t-1} x_s^{D,b} \right. \right. \\ &\quad \left. \left. + \mu_{t-1} - D_{t-1} - x_t^{D,b} \right)^+ \right] \\ &\leq \mathbf{E} \left[\left(\sum_{s=1}^{t-2} \mu_s + \sum_{s=2}^{t-1} x_s^{H_b,b} - \sum_{s=1}^{t-2} D_s - \sum_{s=2}^{t-1} x_s^{D,b} \right)^+ \right. \\ &\quad \left. + (\mu_{t-1} - D_{t-1} - x_t^{D,b})^+ \right] \\ &= \mathbf{E} \left[\left(\sum_{s=1}^{t-2} \mu_s + \sum_{s=2}^{t-2} x_s^{H_b,b} - \sum_{s=1}^{t-2} D_s - \sum_{s=2}^{t-1} x_s^{D,b} \right)^+ \right. \\ &\quad \left. + (\mu_{t-1} - D_{t-1} - x_t^{D,b})^+ \right] \\ &= \mathbf{E} \left[\sum_{s=2}^{t-1} (x_s^{H_b,b} - x_s^{D,b}) + (\mu_{t-1} - D_{t-1} - x_t^{D,b})^+ \right] \\ &\leq \mathbf{E} \left[\sum_{s=1}^{t-1} (\mu_s - x_{s+1}^{D,b} - D_s)^+ \right] \end{aligned}$$

where the third equality follows from the definition of $x_{t-1}^{H_b,b}$ and the identity $(-z + z^+)^+ = (-z)^+$ for all z , the fourth equality follows from the identity in (C.1), and the last inequality follows by induction hypothesis in (C.2). This completes the proof. \blacksquare

Proof of Lemma II.6. For $t \geq 2$, we have:

$$\begin{aligned}
x_t^{H_b,b} &= \left(y_t^{D,b} + \sum_{s=t-L}^t x_s^{D,b} - I_t^{H_b,b} - \sum_{s=t-L}^{t-1} x_s^{H_b,b} \right)^+ \\
&= \left((y_{t-1}^{D,b} + x_{t-1-L}^{D,b} - \mu_{t-1}) + \sum_{s=t-L}^t x_s^{D,b} - (I_{t-1}^{H_b,b} + x_{t-1-L}^{H_b,b} - D_{t-1}) - \sum_{s=t-L}^{t-1} x_s^{H_b,b} \right)^+ \\
&= \left(x_t^{D,b} + \Delta_{t-1} + y_{t-1}^{D,b} - I_{t-1}^{H_b,b} + \sum_{s=t-1-L}^{t-1} x_s^{D,b} - \sum_{s=t-1-L}^{t-1} x_s^{H_b,b} \right)^+ \\
&= \left(x_t^{D,b} + \Delta_{t-1} + U_{t-1}^b \right)^+, \tag{C.3}
\end{aligned}$$

where the second equality follows from the definition of $y_t^{D,b}$ and $I_t^{H_b,b}$, and the last equality follows by the definition of U_{t-1}^b . In addition, we also have:

$$\begin{aligned}
U_t^b &= y_t^{D,b} - I_t^{H_b,b} + \sum_{s=t-L}^t (x_s^{D,b} - x_s^{H_b,b}) \\
&= \left[(y_{t-1}^{D,b} + x_{t-1-L}^{D,b} - \mu_{t-1}) - (I_{t-1}^{H_b,b} + x_{t-1-L}^{H_b,b} - D_{t-1}) + \sum_{s=t-L}^{t-1} (x_s^{D,b} - x_s^{H_b,b}) \right] + x_t^{D,b} - x_t^{H_b,b} \\
&= \left[y_{t-1}^{D,b} - I_{t-1}^{H_b,b} + \sum_{s=t-1-L}^{t-1} (x_s^{D,b} - x_s^{H_b,b}) + D_{t-1} - \mu_{t-1} \right] + x_t^{D,b} - (x_t^{D,b} + \Delta_{t-1} + U_{t-1}^b)^+ \\
&= U_{t-1}^b + \Delta_{t-1} + x_t^{D,b} - (x_t^{D,b} + \Delta_{t-1} + U_{t-1}^b)^+ \\
&= -(U_{t-1}^b + \Delta_{t-1} + x_t^{D,b})^-,
\end{aligned}$$

where the second equality follows from the definition of $y_t^{D,b}$ and $I_t^{H_b,b}$, the third equality follows from (C.3), and the fourth equality follows from the definition of U_{t-1} , and the last equality follows from fact that $-(z)^- = z - (z)^+$ for all z . \blacksquare

Proof of Lemma II.7. For $1 \leq t \leq L+2$, as $x_s^{H_b,b} = x_s^{D,b}$ for $s \leq 1$, it is obvious that $\mathbf{E}[I_t^{H_b,b}] = y_t^{D,b}$. For $L+3 \leq t \leq T$, we prove by induction.

For $t = L + 3$, note that

$$\begin{aligned}
\mathbf{E}[I_{L+3}^{H_b,b}] &= \mathbf{E}[I_{L+2}^{H_b,b} + x_2^{H_b,b} - D_{L+2}] \\
&= \mathbf{E}[I_{L+2}^{H_b,b} + (x_2^{D,b} + \Delta_1 + U_1^b)^+ - D_{L+2}] \\
&= \mathbf{E}[I_{L+2}^{H_b,b} + (x_2^{D,b} + \Delta_1 + U_1^b) - D_{L+2} + (x_2^{D,b} + \Delta_1 + U_1^b)^-] \\
&= y_{L+2}^{D,b} + x_2^{D,b} - \mu_{L+2} + \mathbf{E}[U_1^b - U_2^b] \\
&= y_{L+3}^{D,b} + k_1^b,
\end{aligned}$$

where the second equality follows by the identity of $x_2^{H_b,b}$ in Lemma II.6 and the fourth equality follows by $\mathbf{E}[\Delta_1] = 0$ and the identity of U_2^b in Lemma II.6.

Now, suppose that $\mathbf{E}[I_t^{H_b,b}] = y_t^{D,b} + \sum_{s=2}^{t-1-L} k_{s-1}^b$ for all $t \leq t'$. Then, for $t = t' + 1$, by Lemma II.6 and our induction hypothesis,

$$\begin{aligned}
\mathbf{E}[I_t^{H_b,b}] &= \mathbf{E}[I_{t-1}^{H_b,b} + x_{t-1-L}^{H_b,b} - D_{t-1}] \\
&= \mathbf{E}[I_{t-1}^{H_b,b} + (x_{t-1-L}^{D,b} + \Delta_{t-2-L} + U_{t-2-L}^b)^+ - D_{t-1}] \\
&= \mathbf{E}[I_{t-1}^{H_b,b} + (x_{t-1-L}^{D,b} + \Delta_{t-2-L} + U_{t-2-L}^b) - D_{t-1} + (-x_{t-1-L}^{D,b} - \Delta_{t-2-L} - U_{t-2-L}^b)^+] \\
&= \left[y_{t-1}^{D,b} + \sum_{s=2}^{t-2-L} k_{s-1}^b \right] + x_{t-1-L}^{D,b} - \mu_{t-1} + \mathbf{E}[U_{t-2-L}^b - U_{t-1-L}^b] \\
&= y_{t-1}^{D,b} + \sum_{s=2}^{t-2-L} k_{s-1}^b + x_{t-1-L}^{D,b} - \mu_{t-1} + k_{t-2-L}^b \\
&= y_t^{D,b} + \sum_{s=2}^{t-1-L} k_{s-1}^b,
\end{aligned}$$

where the second equality follows by the identity of $x_{t-1-L}^{H_b,b}$ in Lemma II.6 and the fourth equality follows by $\mathbf{E}[\Delta_{t-2-L}] = 0$ and the identity of U_{t-1-L}^b in Lemma II.6.

This completes the induction. \blacksquare

Proof of Lemma II.8. For notational brevity, we will suppress the notational depen-

dency on α . For $t \leq 1$, we have $x_t^{H_\ell, \ell} = x_t^{D, b}$, so the inequality $P(I_t^{H_\ell, \ell} + x_{t-L}^{H_\ell, \ell} - D_t \geq 0) \geq 1 - \alpha$ naturally holds.

For $t \geq 2$, we can bound:

$$\begin{aligned} P\left(I_t^{H_\ell, \ell} + x_{t-L}^{H_\ell, \ell} - D_t \geq 0 \mid \mathfrak{S}_{t-L-1}^{H_\ell, \ell}\right) &\geq P\left(I_{t-L}^{H_\ell, \ell} + \sum_{s=t-L}^t x_{s-L}^{H_\ell, \ell} - \sum_{s=t-L}^t D_s \geq 0 \mid \mathfrak{S}_{t-L-1}^{H_\ell, \ell}\right) \\ &\geq P\left(y_t^{D, b} + \sum_{s=t-L}^t x_s^{D, b} - \sum_{s=t-L}^t D_s \geq 0 \mid \mathfrak{S}_{t-L-1}^{H_\ell, \ell}\right) \\ &\geq 1 - \alpha \end{aligned}$$

where the first inequality follows from the fact that, given the same starting point and the same order quantities, the inventory level in the lost-sale system is larger than the inventory level in the backorder system, the second inequality follows from the definition of H_ℓ that $x_t^{H_\ell, \ell}(\alpha) = \left(y_t^{D, b}(\alpha) + \sum_{s=t-L}^t x_s^{D, b}(\alpha) - I_t^{H_\ell, \ell} - \sum_{s=t-L}^{t-1} x_s^{H_\ell, \ell}(\alpha)\right)^+$, and the last inequality is directly from the deterministic system. \blacksquare

Proof of Lemma II.11. As the order-up-to policies H_ℓ and H_b both use the same order-up-to level, we have

$$I_t^{H_\ell, \ell} + \sum_{s=t-L}^t x_s^{H_\ell, \ell} = I_t^{H_b, b} + \sum_{s=t-L}^t x_s^{H_b, b}. \quad (\text{C.4})$$

We will simultaneously prove the following set of inequalities by induction:

$$I_t^{H_\ell, \ell} \geq I_t^{H_b, b}, \quad (\text{C.5})$$

$$I_t^{H_\ell, \ell} + \sum_{s=t-L}^{t-1} x_s^{H_\ell, \ell} \geq I_t^{H_b, b} + \sum_{s=t-L}^{t-1} x_s^{H_b, b}, \quad (\text{C.6})$$

$$x_t^{H_\ell, \ell} \leq x_t^{H_b, b}, \quad (\text{C.7})$$

$$I_t^{H_\ell, \ell} + x_{t-L}^{H_\ell, \ell} \geq I_t^{H_b, b} + x_{t-L}^{H_b, b}. \quad (\text{C.8})$$

For our induction base, note that $0 = I_1^{H_\ell, \ell} \geq I_1^{H_b, b} = 0$ and $x_t^{H_\ell, \ell} = x_t^{H_b, b}$ for $1 - L \leq t \leq 1$. So, inequalities (C.5) - (C.8) naturally hold for $t = 1$. Next, suppose that inequalities (C.5) - (C.8) hold for all $t \leq t'$. For $t = t' + 1$, we have the following:

1. For inequality (C.5),

$$\begin{aligned}
I_t^{H_\ell, \ell} &= (I_{t-1}^{H_\ell, \ell} + x_{t-1-L}^{H_\ell, \ell} - D_{t-1})^+ \\
&\geq I_{t-1}^{H_\ell, \ell} + x_{t-1-L}^{H_\ell, \ell} - D_{t-1} \\
&\geq I_{t-1}^{H_b, b} + x_{t-1-L}^{H_b, b} - D_{t-1} \\
&= I_t^{H_b, b},
\end{aligned}$$

where the second inequality follows from the induction hypothesis on inequality (C.8).

2. For inequality (C.6),

$$\begin{aligned}
I_t^{H_\ell, \ell} + \sum_{s=t-L}^{t-1} x_s^{H_\ell, \ell} &= (I_{t-1}^{H_\ell, \ell} + x_{t-1-L}^{H_\ell, \ell} - D_{t-1})^+ + \sum_{s=t-L}^{t-1} x_s^{H_\ell, \ell} \\
&\geq I_{t-1}^{H_\ell, \ell} + x_{t-1-L}^{H_\ell, \ell} - D_{t-1} + \sum_{s=t-L}^{t-1} x_s^{H_\ell, \ell} \\
&= I_{t-1}^{H_b, b} + x_{t-1-L}^{H_b, b} - D_{t-1} + \sum_{s=t-L}^{t-1} x_s^{H_b, b},
\end{aligned}$$

where the second equality is from (C.4).

3. Inequality (C.7) follows immediately from (C.4) and (C.6).

4. For inequality (C.8), first note that, by induction hypothesis, $x_s^{H_\ell, \ell} \leq x_s^{H_b, b}$ for all

$s \leq t'$. Since we also have $x_t^{H_\ell, \ell} \leq x_t^{H_b, b}$ (as argued in #3), this implies:

$$\sum_{s=t-L+1}^t x_s^{H_\ell, \ell} \leq \sum_{s=t-L+1}^t x_s^{H_b, b}.$$

Put the above inequality together with (C.4) immediately yields $I_t^{H_\ell, \ell} + x_{t-L}^{H_\ell, \ell} \geq I_t^{H_b, b} + x_{t-L}^{H_b, b}$.

This completes our induction. \blacksquare

Proof of Lemma II.13. The proof is very similar to the proof of Lemma II.2; we omit the details. \blacksquare

Proof of Lemma II.14. We prove by induction. For $t = 1$, $I_1^{\pi, \ell} + x_{1-L}^{\pi, \ell} - D_1 = -\Upsilon_1$. Suppose that the identity in Lemma II.14 is true for any $t \leq t'$. Then, for $t = t' + 1$, we have:

$$\begin{aligned} I_t^{\pi, \ell} + x_{t-L}^{\pi, \ell} - D_t &= (I_{t-1}^{\pi, \ell} + x_{t-1-L}^{\pi, \ell} - D_{t-1})^+ + x_{t-L}^{\pi, \ell} - D_t \\ &= \max \left\{ 0, -\sum_{s=1}^{t-1} \Upsilon_s + \max \left\{ 0, \Upsilon_1, \sum_{s=1}^2 \Upsilon_s, \dots, \sum_{s=1}^{t-2} \Upsilon_s \right\} \right\} - \Upsilon_t \\ &= -\sum_{s=1}^t \Upsilon_s + \max \left\{ 0, \Upsilon_1, \sum_{s=1}^2 \Upsilon_s, \dots, \sum_{s=1}^{t-2} \Upsilon_s, \sum_{s=1}^{t-1} \Upsilon_s \right\}, \end{aligned}$$

where the first equation follows by induction hypothesis. This completes the induction. \blacksquare

Proof of Lemma II.15. Fix $\pi \in \Pi^\ell$. We want to show that π satisfies the probabilistic service level constraints defined in $\Pi^{\tilde{b}}$ for all $1 \leq t \leq T$ (i.e., if we apply the same ordering decision x_t^π at period t in the backorder system \tilde{b} as if the inventory levels evolve according to a lost-sales system, the resulting sequence of ordering decisions satisfies the probabilistic service level constraints in $\Pi^{\tilde{b}}$). In what follows, we will

consider the two cases $1 \leq t \leq L + 1$ and $L + 2 \leq t \leq T$ in turn; their proofs proceed by induction.

Step 1: Case $1 \leq t \leq L + 1$

For $t = 1$, since $P(x_{1-L}^\pi - D_1 \geq 0) \geq 1 - \alpha$, the probabilistic service level constraint in $\Pi^{\tilde{b}}$ is also satisfied. Suppose that the probabilistic service level constraints in $\Pi^{\tilde{b}}$ are satisfied for $1 \leq t \leq t' < L + 1$, i.e.,

$$1 - \alpha \leq P\left(I_t^{\pi, \tilde{b}} + x_{t-L}^\pi - D_t \geq -\sum_{s=1}^{t-1} \theta_s(\alpha)\right) = P\left(\sum_{s=1}^t x_{s-L}^\pi - \sum_{s=1}^t D_s \geq -\sum_{s=1}^{t-1} \theta_s(\alpha)\right)$$

for $1 \leq t \leq t' < L + 1$. By definition of $\gamma_t(\alpha)$, this implies:

$$\sum_{s=1}^t x_{s-L}^\pi + \sum_{s=1}^{t-1} \theta_s \geq \sum_{s=1}^t \mu_s + \beta_1^t(\alpha) \quad (\text{C.9})$$

for $1 \leq t \leq t' < L + 1$. So, for $t = t' + 1$, we can bound:

$$\begin{aligned} & P\left(I_t^{\pi, \tilde{b}} + x_{t-L}^\pi - D_t \geq -\sum_{s=1}^{t-1} \theta_s(\alpha)\right) \\ &= P\left(\sum_{s=1}^t x_{s-L}^\pi - \sum_{s=1}^t D_s \geq -\sum_{s=1}^{t-1} \theta_s(\alpha)\right) \\ &= P\left(\sum_{s=1}^t x_{s-L}^\pi - \sum_{s=1}^t D_s \geq -\max\left\{0, \theta_1(\alpha), \dots, \sum_{s=1}^{t-1} \theta_s(\alpha)\right\}\right) \\ &\geq P\left(\sum_{s=1}^t x_{s-L}^\pi - \sum_{s=1}^t D_s \geq -\max\left\{0, \Upsilon_1, \sum_{s=1}^2 \Upsilon_s, \dots, \sum_{s=1}^{t-1} \Upsilon_s\right\}\right) \\ &= P(I_t^{\pi, \ell} + x_{t-L}^\pi - D_t \geq 0) \\ &\geq 1 - \alpha, \end{aligned}$$

where the second equality follows since $\theta_s(\alpha) \geq 0$; the first inequality follows by equa-

tion (C.9) and

$$\sum_{k=1}^s \Upsilon_k = \sum_{k=1}^s D_k - \sum_{k=1}^s x_{k-L}^\pi \leq s\bar{D} - \sum_{k=1}^s \mu_k - \beta_1^s(\alpha) + \sum_{k=1}^{s-1} \theta_k(\alpha) = \sum_{k=1}^s \theta_k(\alpha)$$

the third equality follows by Lemma II.14; and, the last inequality follows since $\pi \in \Pi^\ell$. This completes the induction. We conclude that the probabilistic service level constraints in $\Pi^{\bar{b}}$ are satisfied by π for $1 \leq t \leq L+1$.

Step 2: Case $L+2 \leq t \leq T$

For $t = L+2$, we have:

$$\begin{aligned} & P \left(I_{L+2}^{\pi, \bar{b}} + x_2^\pi - D_{L+2} \geq - \sum_{s=1}^{L+1} \theta_s(\alpha) \mid \mathfrak{F}_1^{\pi, \bar{b}} \right) \\ &= P \left(\sum_{s=1}^{L+2} x_{s-L}^\pi - \sum_{s=1}^{L+2} D_s + \max \left\{ 0, \theta_1(\alpha), \dots, \sum_{s=1}^{L+1} \theta_s(\alpha) \right\} \geq 0 \right) \\ &\geq P \left(\sum_{s=1}^{L+2} x_{s-L}^\pi - \sum_{s=1}^{L+2} D_s + \max \left\{ 0, \Upsilon_1, \sum_{s=1}^2 \Upsilon_s, \dots, \sum_{s=1}^{L+1} \Upsilon_s \right\} \geq 0 \right) \\ &= P(I_t^{\pi, \ell} + x_{t-L}^\pi - D_t \geq 0) \\ &\geq 1 - \alpha, \end{aligned}$$

where the first equality holds since $\theta_s(\alpha) \geq 0$; the first inequality again follows from (C.9), which holds for $t = L+1$ (by our result in Step 1 above); the second equality follows from Lemma II.14; and, the last inequality follows since $\pi \in \Pi^\ell$.

Suppose that the second probabilistic service level constraints in $\Pi^{\bar{b}}$ are satisfied

for $L + 2 \leq t \leq t' < T$; in particular,

$$\begin{aligned} 1 - \alpha &\leq P \left(I_t^{\pi, \tilde{b}} + x_{t-L}^{\pi} - D_t \geq - \sum_{s=1}^{t-1} \theta_s(\alpha) \mid \mathfrak{S}_{t-L-1}^{\pi, \tilde{b}} \right) \\ &= P \left(I_{t-L}^{\pi, \tilde{b}} + \sum_{s=t-L}^t x_{s-L}^{\pi} - \sum_{s=t-L}^t D_s \geq - \sum_{s=1}^{t-1} \theta_s(\alpha) \mid \mathfrak{S}_{t-L-1}^{\pi} \right) \end{aligned}$$

for $L + 2 \leq t \leq t' < T$. By definition of $w_t(\alpha)$, this implies:

$$I_{t-L}^{\pi, \tilde{b}} + \sum_{s=t-L}^t x_{s-L}^{\pi} + \sum_{s=1}^{t-1} \theta_s(\alpha) \geq \sum_{s=t-L}^t \mu_s + \beta_{t-L}^{L+1}(\alpha)$$

for $L + 2 \leq t \leq t' < T$. So, for $t = t' + 1$, we can bound:

$$\begin{aligned} &P \left(I_t^{\pi, \tilde{b}} + x_{t-L}^{\pi} - D_t \geq - \sum_{s=1}^{t-1} \theta_s(\alpha) \mid \mathfrak{S}_{t-L-1}^{\pi, \tilde{b}} \right) \\ &= P \left(I_t^{\pi, \tilde{b}} + x_{t-L}^{\pi} - D_t \geq - \max \left\{ 0, \theta_1(\alpha), \sum_{s=1}^2 \theta_s(\alpha), \dots, \sum_{s=1}^{t-1} \theta_s(\alpha) \right\} \mid \mathfrak{S}_{t-L-1}^{\pi, \tilde{b}} \right) \\ &\geq P \left(I_t^{\pi, \tilde{b}} + x_{t-L}^{\pi} - D_t \geq - \max \left\{ 0, \Upsilon_1, \sum_{s=1}^2 \Upsilon_s, \dots, \sum_{s=1}^{t-1} \Upsilon_s \right\} \mid \mathfrak{S}_{t-L-1}^{\pi, \ell} \right) \\ &= P \left(I_t^{\pi, \ell} + x_{t-L}^{\pi} - D_t \geq 0 \mid \mathfrak{S}_{t-L-1}^{\pi, \ell} \right) \\ &\geq 1 - \alpha, \end{aligned}$$

where the first equality follows since $\theta_s(\alpha) \geq 0$; the first inequality follows from (C.10),

i.e.,

$$\sum_{k=1}^s \Upsilon_k = \sum_{k=s-L}^s (D_k - x_{k-L}^{\pi}) - I_{s-L}^{\pi, \tilde{b}} \leq (L+1)\bar{D} - \sum_{k=s-L}^s \mu_k - \beta_{s-L}^{L+1}(\alpha) + \sum_{k=1}^{s-1} \theta_k(\alpha) = \sum_{k=1}^s \theta_k(\alpha);$$

the third equality follows from Lemma II.14; and, the last inequality follows since $\pi \in \Pi^{\ell}$. This completes the induction; we conclude that the probabilistic service level constraints in $\Pi^{\tilde{b}}$ are also satisfied by π for $L + 2 \leq t \leq T$, which proves Lemma II.15.

■

Proof of Lemma II.16. Denote the optimal control corresponding to $C^{*,\ell}$ as π^* and its order quantities as $\{x_t^{\pi^*}\}_t$. Now let us construct a new control $\tilde{\pi}$ such that

$$x_{t+1}^{\tilde{\pi}} \triangleq x_{t+1}^{\pi^*} + \left(D_t - x_{t-L}^{\pi^*} - I_t^{\pi^*,\ell} \right)^+,$$

namely, the order quantity under $\tilde{\pi}$ equals the order quantity under π^* plus the amount of lost-sales incurred during the previous period in the optimal lost-sales system. As $\tilde{\pi}$ orders at least as many as π^* does and π^* is feasible to the backorder system \tilde{b} by Lemma II.15, $\tilde{\pi}$ is also feasible to the backorder system \tilde{b} . It follows that

$$C^{*,\tilde{b}} - C^{*,\ell} \leq C^{\tilde{\pi},\tilde{b}} - C^{*,\ell}.$$

Let us compare the backorder system \tilde{b} under policy $\tilde{\pi}$ (called System I) to the lost-sales system under the optimal policy π^* (called System II). Observe that

1. the two systems always have the same inventory position;
2. System I always has a lower inventory level than System II does;
3. the additional total penalty costs in System I (compared with System II) is at most pL times the total lost-sales incurred in System II.

Note that a similar comparison between a backlog system and a lost-sales one was

conducted earlier in the cost-based model (see Theorem 5 in 48). Hence, we have:

$$\begin{aligned}
C^{\bar{\pi}, \bar{b}} - C^{*, \ell} &\leq c \cdot \sum_{t=1}^T \mathbb{E} \left[\left(D_t - x_{t-L}^{\pi^*} - I_t^{\pi^*, \ell} \right)^+ \right] + L \cdot p \cdot \sum_{t=1}^T \mathbb{E} \left[\left(D_t - x_{t-L}^{\pi^*} - I_t^{\pi^*, \ell} \right)^+ \right] \\
&= (c + L \cdot p) \cdot \sum_{t=1}^T \mathbb{E} \left[\mathbb{I} \left(D_t - x_{t-L}^{\pi^*} - I_t^{\pi^*, \ell} \geq 0 \right) \left(D_t - x_{t-L}^{\pi^*} - I_t^{\pi^*, \ell} \right) \right] \\
&\leq (c + L \cdot p) \cdot T \cdot \alpha \cdot \bar{D},
\end{aligned}$$

where the first inequality comes from the observations above and the last inequality comes from the feasibility of the policy π^* . The proof is completed. \blacksquare

Proof of Lemma II.17. We first re-write **DET** in matrix form. Replace y_t in **DET** with $\sum_{s=1}^{t-1} x_{s-L} - \sum_{s=1}^{t-1} \mu_s$ (this way, we no longer need the first two constraints in (2.5)). Now, (2.5) can be expressed as follows:

$$\begin{aligned}
D^{*,b}(\alpha) &= \min_{x,z,m} \begin{bmatrix} d_1 \\ d_2 \\ d_3 \end{bmatrix}' \begin{bmatrix} x \\ z \\ m \end{bmatrix} \\
\text{s.t.} & \begin{bmatrix} A & O & O \\ -A & I & O \\ A & O & I \end{bmatrix} \begin{bmatrix} x \\ z \\ m \end{bmatrix} \geq \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}, \quad \begin{bmatrix} x \\ z \\ m \end{bmatrix} \geq 0,
\end{aligned}$$

where d_1 is a vector of c 's, d_2 is a vector of h 's, d_3 is a vector of p 's, the b 's are the proper right-hand-side (RHS) of the constraints in (2.5), I is an identity matrix, O is a zero matrix.

We can also rewrite formulation (2.16) compactly as follows:

$$D^{*,\tilde{b}} = \min_{x,z,m} \begin{bmatrix} d_1 \\ d_2 \\ d_3 \end{bmatrix}' \begin{bmatrix} x \\ z \\ m \end{bmatrix}$$

$$\text{s.t.} \quad \begin{bmatrix} A & O & O \\ -A & I & O \\ A & O & I \end{bmatrix} \begin{bmatrix} x \\ z \\ m \end{bmatrix} \geq \begin{bmatrix} \tilde{b}_1 \\ \tilde{b}_2 \\ \tilde{b}_3 \end{bmatrix}, \quad \begin{bmatrix} x \\ z \\ m \end{bmatrix} \geq 0,$$

where $\tilde{b}_2 = b_2$ and $\tilde{b}_3 = b_3$. The dual formulation of (2.16) is given by

$$\max_{\lambda_1, \lambda_2, \lambda_3} \begin{bmatrix} \tilde{b}_1 \\ \tilde{b}_2 \\ \tilde{b}_3 \end{bmatrix}' \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \lambda_3 \end{bmatrix} \quad (:= \tilde{b}'\lambda)$$

$$\text{s.t.} \quad \begin{bmatrix} A' & -A' & A' \\ O & I & O \\ O & O & I \end{bmatrix} \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \lambda_3 \end{bmatrix} \leq \begin{bmatrix} d_1 \\ d_2 \\ d_3 \end{bmatrix}, \quad \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \lambda_3 \end{bmatrix} \geq 0.$$

Let λ^D and $\tilde{\lambda}^D$ denote optimal dual solutions to $D^{*,b}$ and $D^{*,\tilde{b}}$, respectively. By the standard duality theory,

$$D^{*,b} - D^{*,\tilde{b}} = b'\lambda^D - \tilde{b}'\tilde{\lambda}^D \leq (b - \tilde{b})'\lambda^D = (b_1 - \tilde{b}_1)'\lambda_1^D,$$

where the inequality follows because, by the optimality of $\tilde{\lambda}^D$, we have $\tilde{b}'\tilde{\lambda}^D \geq \tilde{b}'\lambda^D$, and the last equality follows since $\tilde{b}_2 = b_2$ and $\tilde{b}_3 = b_3$. Let $\bar{\Delta} = b_1 - \tilde{b}_1$. We can further

bound $D^{*,b} - D^{*,\bar{b}}$ as follows:

$$\begin{aligned}
D^{*,b} - D^{*,\bar{b}} &\leq \max_{q_1, q_2, q_3} \bar{\Delta}' q_1 \\
&\text{s.t. } \begin{bmatrix} A' & -A' & A' \\ O & I & O \\ O & O & I \end{bmatrix} \begin{bmatrix} q_1 \\ q_2 \\ q_3 \end{bmatrix} \leq \begin{bmatrix} d_1 \\ d_2 \\ d_3 \end{bmatrix}, \quad \begin{bmatrix} q_1 \\ q_2 \\ q_3 \end{bmatrix} \geq 0 \\
&\leq \max_{q_1} \bar{\Delta}' q_1 \tag{C.10} \\
&\text{s.t. } A' q_1 \leq A' d_2 + d_1, \quad q_1 \geq 0.
\end{aligned}$$

The second inequality above follows by setting $q_2 = d_2$ and $q_3 = 0$. Now, note that

$$A' d_2 + d_1 = \begin{bmatrix} c + T \cdot h \\ c + (T-1) \cdot h \\ \vdots \\ c + h \end{bmatrix} \quad \text{and} \quad \bar{\Delta} = \begin{bmatrix} 0 \\ \theta_1(\alpha) \\ \vdots \\ \sum_{s=1}^{T-1} \theta_s(\alpha) \end{bmatrix}.$$

Since $\theta_t(\alpha) \geq 0$, it is not difficult to see that the optimal solution to (C.10) is given by

$$q_1 = \begin{bmatrix} h \\ h \\ \vdots \\ h \\ c + h \end{bmatrix}.$$

We conclude that

$$D^{*,b} - D^{*,\bar{b}} \leq (c + h) \cdot \left[\sum_{t=1}^{T-1} \theta_t(\alpha) \right] + h \cdot \left[\sum_{t=1}^{T-1} \sum_{s=1}^{t-1} \theta_s(\alpha) \right] \leq (c + h) \cdot \phi(T, \alpha). \quad \blacksquare$$

APPENDIX D

Proof of Chapter III

Proof of Section 3.4

Proof of Proposition III.1:

In this case, any pricing decision corresponds a certain demand structure – (1) The snobs purchase product 1 while the followers purchase product 2. (2) Both snobs and followers purchase product 1. (3) Both snobs and followers purchase product 2. (4) Snobs purchase 1 while followers make no purchase. (5) Snobs purchase 2 while followers make no purchase.

Scenario (1). The optimization problem that the firm faces is $\max_{p_1, p_2 \geq 0} p_1 + \beta p_2$, s.t. $p_1 - p_2 \leq v_S \Delta q$, $p_1 \leq v_S q_1$, $p_1 - p_2 \geq v_F \Delta q$, $p_2 \leq v_F q_2$. The optimal prices are $p_1^* = v_S q_1 - (v_S - v_F) q_2$ and $p_2^* = v_F q_2$. The optimal revenue is $R^* = v_S \Delta q + (1 + \beta) v_F q_2$.

Scenario (2). The optimization problem that the firm faces is $\max_{p_1, p_2 \geq 0} (1 + \beta) p_1$, s.t. $p_1 - p_2 \leq v_S \Delta q$, $p_1 \leq v_S q_1$, $p_1 - p_2 \leq v_F \Delta q$, $p_1 \leq v_F q_1$. The optimal prices are $p_1^* = v_F q_1$ and $p_2^* > v_F q_2$. The optimal revenue is $R^* = (1 + \beta) v_F q_1$.

Scenario (3). The optimization problem that the firm faces is $\max_{p_1, p_2 \geq 0} (1 + \beta) p_2$, s.t. $p_1 - p_2 \geq v_S \Delta q$, $p_2 \leq v_S q_1$, $p_1 - p_2 \geq v_F \Delta q$, $p_2 \leq v_F q_1$. The optimal prices are $p_1^* > v_F q_2 + v_S \Delta q$ and $p_2^* = v_F q_2$. The optimal revenue is $R^* = (1 + \beta) v_F q_2$, which is clearly dominated by scenario (2).

Scenario (4). The optimization problem that the firm faces is $\max_{p_1, p_2 \geq 0} p_1$, s.t. $p_1 - p_2 \leq v_S \Delta q$, $p_1 \geq v_F q_1$, $p_1 \geq v_F q_2$, $p_1 \leq v_S q_1$. The optimal prices are $p_1^* = v_S q_1$ and $p_2^* > v_S q_2$. The optimal revenue is $R^* = v_S q_1$.

Scenario (5). The optimization problem that the firm faces is $\max_{p_1, p_2 \geq 0} p_2$, s.t.

$p_1 - p_2 \geq v_S \Delta q$, $p_2 \geq v_F q_2$, $p_1 \geq v_S q_1$, $p_2 \geq v_F q_2$. The optimal prices are $p_1^* \geq v_S q_2$ and $p_2^* = v_S q_2$. The optimal revenue is $R^* = v_S q_2$, which is dominated by scenario 4.

In summary, the optimal policies can only be among scenarios (1), (2) and (4) and we can further reduce it as follows. Note that (a) if $v_S \geq (1 + \beta)v_F$, then the profit in scenario (4) dominate those in (1) and (2), i.e., $v_S q_1 \geq v_S \Delta q + (1 + \beta)v_F q_2$ and $v_S q_1 \geq (1 + \beta)v_F q_1$; (b) if $v_S < (1 + \beta)v_F$, then the profit in scenario (2) dominate those in (1) and (4), i.e., $(1 + \beta)v_F q_1 \geq v_S \Delta q + (1 + \beta)v_F q_2$ and $(1 + \beta)v_F q_1 \geq v_S q_1$.

Proof of Theorem III.2:

We follow the logic of rational expectation in this proof: we first assume a certain expectation of demand, then solve for the corresponding optimal prices such that the realized demand meet the expected demand.

In this setting, there are 8 scenarios of the demand expectation: (1) The snob customers purchase product 1 while the followers purchase product 2. (2) Both snob and follower customers purchase product 1. (3) Both snob and follower customers purchase product 2. (4) Snob customers purchase 1 while follower customers make no purchase. (5) Snob customers purchase 2 while follower customers make no purchase. (6) The snob customers purchase product 2 while the followers purchase product 1. (7) The snob customers make no purchase while the followers purchase product 1. (8) The snob customers make no purchase while the followers purchase product 2. We analyze them individually. In the following analysis, we denote $q_1 - q_2$ simply as Δq .

Scenario (1). The snob customers purchase product 1 while the followers purchase product 2. Thus, $D_1^e = 1$ and $D_2^e = \beta$. The conditions for snob to purchase 1 are $v_S \geq \frac{p_1 - p_2 + k(1 - \beta)}{\Delta q}$ and $v_S \geq \frac{p_1 + k}{q_1}$, while the conditions for the follower customers choose to purchase product 2 are $v_F < \frac{p_1 - p_2 - m(1 - \beta)}{\Delta q}$ and $v_F \geq \frac{p_2 - m\beta}{q_2}$. In other words, the optimization problem that the firm faces is $\max_{p_1, p_2 \geq 0} p_1 + \beta p_2$, s.t. $p_1 - p_2 \leq v_S \Delta q - k(1 - \beta)$, $p_1 \leq v_S q_1 - k$, $p_1 - p_2 \geq v_F \Delta q + m(1 - \beta)$, $p_2 \leq v_F q_2 + m\beta$. The condition for feasibility of the above problem is $(v_S - v_F)\Delta q \geq (k + m)(1 - \beta)$. It is easy to observe that the optimal prices are as follows.

If $v_S q_2 - k\beta \leq v_F q_2 + m\beta$, $p_1^* = v_S q_1 - k$ and $p_2^* = \min\{v_F q_2 + m\beta, v_S q_1 - k - v_F \Delta q - m(1 - \beta)\}$. O.w., $p_1^* = v_F q_2 + m\beta + v_S \Delta q - k(1 - \beta)$ and $p_2^* = v_F q_2 + m\beta$.

We can write the optimal revenue compactly as $R^* = \min\{v_S q_1 - k, v_F q_2 + m\beta + v_S \Delta q - k(1 - \beta)\} + \beta \min\{v_F q_2 + m\beta, v_S q_1 - k - v_F \Delta q - m(1 - \beta)\}$.

Scenario (2). Both snob and follower customer purchase product 1. $D_1^e = 1 + \beta$ and $D_2^e = 0$.

The optimization problem that the firm faces is $\max_{p_1, p_2 \geq 0} (1 + \beta)p_1$, s.t. $p_1 -$

$p_2 \leq v_S \Delta q - k(1 + \beta)$, $p_1 \leq v_S q_1 - k(1 + \beta)$, $p_1 - p_2 \leq v_F \Delta q + m(1 + \beta)$, $p_1 \leq v_F q_1 + m(1 + \beta)$. The optimal prices are, $p_1^* = \min\{v_S q_1 - k(1 + \beta), v_F q_1 + m(1 + \beta)\}$, and any $p_2 \geq p_1^* - \min\{v_S \Delta q - k(1 + \beta), v_F \Delta q + m(1 + \beta)\}$. The optimal revenue is $R^* = (1 + \beta) \min\{v_S q_1 - k(1 + \beta), v_F q_1 + m(1 + \beta)\}$.

Scenario (3). Both snob and follower customers purchase product 2. $D_1^e = 0$ and $D_2^e = 1 + \beta$. The optimization problem is $\max_{p_1, p_2 \geq 0} (1 + \beta)p_2$, s.t. $p_1 - p_2 > v_S \Delta q + k(1 + \beta)$, $p_2 \leq v_S q_2 - k(1 + \beta)$, $p_1 - p_2 > v_F \Delta q - m(1 + \beta)$, $p_2 \leq v_F q_2 + m(1 + \beta)$. The optimal prices are, $p_2^* = \min\{v_S q_2 - k(1 + \beta), v_F q_2 + m(1 + \beta)\}$ and any $p_1^* \geq p_2^* + \max\{v_S \Delta q + k(1 + \beta), v_F \Delta q - m(1 + \beta)\}$. The optimal revenue is $R^* = (1 + \beta) \min\{v_S q_2 - k(1 + \beta), v_F q_2 + m(1 + \beta)\}$. It is easy to see that such a selling strategy generates less revenue than scenario (2).

Scenario (4). Snob customers purchase 1 while follower customers make no purchase. $D_1^e = 1$ and $D_2^e = 0$. The optimization problem is $\max_{p_1, p_2 \geq 0} p_1$, s.t. $p_1 - p_2 \leq v_S \Delta q - k$, $p_1 \leq v_S q_1 - k$, $p_1 > v_F q_1 + m$, $p_2 > v_F q_2$. The condition for feasibility of the above problem is $(v_S - v_F)q_1 \geq (k + m)$.

The optimal prices are, $p_1^* = v_S q_1 - k$ and any $p_2^* \geq v_S q_2$. The optimal revenue is $R^* = v_S q_1 - k$. Then if $v_S q_2 - k\beta \leq v_F q_2 + m\beta$, then it is dominated by the selling strategy described in Scenario (1).

Scenario (5). Snob customers purchase 2 while follower customers make no purchase. $D_1^e = 0$ and $D_2^e = 1$. The optimization problem is $\max_{p_1, p_2 \geq 0} p_2$, s.t. $p_1 - p_2 > v_S \Delta q + k$, $p_2 \leq v_S q_2 - k$, $p_1 > v_F q_1$, $p_2 > v_F q_2 + m$. The condition for feasibility of the above problem is $(v_S - v_F)q_1 \geq (k + m)$.

The optimal prices are, $p_2^* = v_S q_2 - k$ and any $p_1^* \geq v_S q_1$. The optimal revenue is $R^* = v_S q_2 - k$, which is dominated by the selling strategy described in Scenario (4).

Scenario (6). Snob customers purchase 2 while follower customers purchase product 1. $D_1^e = \beta$ and $D_2^e = 1$. The optimization problem is $\max_{p_1, p_2 \geq 0} \beta p_1 + p_2$, s.t. $p_1 - p_2 > v_S \Delta q - k(\beta - 1)$, $p_2 \leq v_S q_2 - k$, $p_1 - p_2 \leq v_F \Delta q + m(\beta - 1)$, $p_1 \leq v_F q_1 + m\beta$. The condition for feasibility of the above problem is $(v_S - v_F)\Delta q < (k + m)(\beta - 1)$. Note that in the cases where $\beta < 1$, such an inequality can never hold.

The optimal prices are, $p_2^* = v_S q_2 - k$ and $p_1^* = \min\{v_S q_2 - k + v_F \Delta q + m(\beta - 1), v_F q_1 + m\beta\}$. The optimal revenue is $R^* = v_S q_2 - k + \beta \min\{v_S q_2 - k + v_F \Delta q + m(\beta - 1), v_F q_1 + m\beta\}$.

Scenario (7). Snob customers make no purchase while follower customers purchase product 1. $D_1^e = \beta$ and $D_2^e = 0$. The optimization problem is $\max_{p_1, p_2 \geq 0} \beta p_1$, s.t. $p_1 > v_S q_1 - k\beta$, $p_2 > v_S q_2$, $p_1 - p_2 \leq v_F \Delta q + m\beta$, $p_1 \leq v_F q_1 + m\beta$. The condition for

feasibility of the above problem is $(v_S - v_F)q_1 < (k + m)\beta$. The optimal prices are, $p_1^* = v_Fq_1 + m\beta$ and any $p_2^* \geq v_Fq_2$. The optimal revenue is $R^* = \beta(v_Fq_1 + m\beta)$.

Scenario (8). Snob customers make no purchase while follower customers purchase product 2. $D_1^e = 0$ and $D_2^e = \beta$. The optimization problem is $\max_{p_1, p_2 \geq 0} \beta p_2$, s.t. $p_1 > v_Sq_1$, $p_2 > v_Sq_2 - k\beta$, $p_1 - p_2 > v_F\Delta q - m\beta$, $p_2 \leq v_Fq_2 + m\beta$. The condition for feasibility of the above problem is $(v_S - v_F)q_2 < (k + m)\beta$.

The optimal prices are, $p_2^* = v_Fq_2 + m\beta$ and any $p_1^* \geq v_Fq_1$. The optimal revenue is $R^* = \beta(v_Fq_2 + m\beta)$, which is dominated by the revenue in Scenario (7).

After knowing the prices in each of the 8 scenarios, we know that the equilibrium can only be among scenarios (1), (2), (4), (6), and (7). Next, we classify the possible forms of equilibria in four cases, depending on the relation between $(v_S - v_F)q_2$ and $(m + k)\beta$. Case 1: When $(v_S - v_F)q_2 \leq (m + k)\beta$, any scenarios other than (1), (2), (6) and (7), are dominated by other scenarios, as scenario (4) is dominated by (1). Case 2: When $(v_S - v_F)q_2 \geq (m + k)\beta$, any scenarios other than (1), (2), (4), and (6), are dominated by other scenarios, as scenario (7) is dominated by (6).

Proof of Proposition III.4 In the case of $(v_S - v_F)q_2 \leq (m + k)\beta$, with Assumption III.3, we know that $(v_S - v_F)(q_1 - q_2) \geq k + m$. Thus, the feasibility condition for scenario (6) does not hold. Also, scenario (7) is dominated by (2) under Assumption III.3.

Proof of Section 3.5

Proof of Proposition III.5. We first analyze the strategy of firm 1 given the price of firm 2, p_2 . For firm 1, the price to get the snob customers is $p_1 \leq v_S(q_1 - q_2) + p_2$ and $p_1 \leq v_Sq_1$ (such that snobs have positive utility from product 1 and it is higher than that from product 2, i.e., $q_1v_S - p_1 \geq q_2v_S - p_2$) and the price to get the followers is $p_1 \leq v_F(q_1 - q_2) + p_2$ and $p_1 \leq v_Fq_1$ (such that followers have positive utility from product 1 and it is higher than that from product 2, i.e., $q_1v_F - p_1 \geq q_2v_F - p_2$).

Similarly, we can get the pricing and revenue function of firm 2. For firm 2 the price to get the snob customers is $p_2 \leq p_1 - v_S(q_1 - q_2)$ and $p_2 \leq v_Sq_2$, and the price to get the followers is $p_2 \leq p_1 - v_F(q_1 - q_2)$ and $p_2 \leq v_Fq_2$.

From the above utility functions, we can observe that the equilibrium may be (and can only be) in the following two forms. Form (1), firm 1 sells to snob customer while firm 2 sells to follower customers. In other words, p_1 has the following form:

$$p_1 = \begin{cases} v_S(q_1 - q_2) + p_2 & \text{if } p_2 \leq v_Sq_2 \\ v_Sq_1 & \text{o.w.} \end{cases}$$

p_2 has the following form:

$$p_2 = \begin{cases} p_1 - v_F(q_1 - q_2) & \text{if } p_1 \leq v_F q_1 \\ v_F q_2 & \text{o.w.} \end{cases}$$

Thus, in equilibrium, firm 1 sets $p_1 = v_S(q_1 - q_2) + v_F q_2$ and firm 2 sets $p_2 = v_F q_2$. The revenue of the firms are $R_1 = v_S(q_1 - q_2) + v_F q_2$ and $R_2 = \beta v_F q_2$. A quick check for the equilibrium is that: For firm 1, rising p_1 means losing snobs customer while lowering p_1 means losing revenue, unless lowering the price all the way to the form (2) shown below; For firm 2, rising p_2 means losing followers while lowering p_2 means losing revenue.

Form (2), firm 1 sets $p_1 = v_F(q_1 - q_2)$ and $p_2 = 0$. Remember that customers choose product 1 if they have equal utility from both products. The revenue of the firms are $R_1 = (1 + \beta)v_F(q_1 - q_2)$ and $R_2 = 0$. A quick check for the equilibrium: For firm 1, rising p_1 means losing customers while lowering p_1 means losing revenue; For firm 2, rising p_2 does not win back any customers while there is no room to further lowering p_2 .

In summary, if $(v_S - v_F)(q_1 - q_2) \geq v_F(\beta q_1 - (1 + \beta)q_2)$, then firm 1 will choose to adopt the form (1) equilibrium (which the firms 2 knows that firm 1 will choose such an equilibrium); Otherwise, firm 1 will choose to adopt the form (2) equilibrium.

Proof of Lemma III.6 We first show that for any given $p_1 > 0$, a pricing war cannot survive in an equilibrium. We prove it by contradiction. Suppose there exist a case where for $p_2 = 0$, firm 1 still occupies the whole market with some strictly positive price $p_1 > 0$. In this case, for any type t ($t \geq 0$) of customers, he has a positive utility from product 1, $tq_1 - p_1 \geq 0$. In such a case, firm 2 can increase p_2 to $p_\epsilon (> 0)$ to attract customers of type v where $0 \leq v \leq \frac{p_1 - p_\epsilon}{\Delta q}$. Any customer with type v prefers product 2 over product 1 and extracts a positive utility from product 2, since $vq_2 - p_\epsilon \geq vq_1 - p_1 \geq 0$ (where the last inequality follows from the condition $tq_1 - p_1 \geq 0$ for any $t \geq 0$). Thus, firm 2 can always deviate from $p_2 = 0$ and get a positive revenue. Next, it is easy to see that if firm 1 sets $p_1 = 0$, he can take the entire market but only with zero revenue. Combining these two parts, we complete the proof.

Proof of Proposition III.7

Remember that for a customer with type v , either snob or followers, he purchases product 1 from firm 1 if the following conditions hold: (1) $v \geq \frac{p_1 - p_2}{\Delta q}$ (such that $q_1 v - p_1 \geq q_2 v - p_2$) and (2) $v \geq \frac{p_1}{q_1}$. Similarly, he purchase product 2 from firm 2 if

the following conditions hold: (1) $v < \frac{p_1 - p_2}{\Delta q}$ and (2) $v \geq \frac{p_2}{q_2}$.

Based on the customer purchase decision, between product 1 and 2, we described in Section 3.4, we can write the aggregation demand in the following two cases:

(1) if $\frac{q_1}{p_1} \leq \frac{q_2}{p_2}$ (product 1 has higher unit cost of quality), then customers who have type $v \geq \frac{p_1 - p_2}{\Delta q}$ purchase product 1, and customers who have type $\frac{p_2}{q_2} \leq v \leq \frac{p_1 - p_2}{\Delta q}$ purchase product 2. The total purchase from follower customers, D_1^F and D_2^F , and those from snob customers, D_1^S and D_2^S , are as follows, respectively.

$$\begin{aligned} D_1^F &= \beta \left(1 - \frac{p_1 - p_2}{\Delta q}\right) \\ D_2^F &= \beta \left(\frac{p_1 - p_2}{\Delta q} - \frac{p_2}{q_2}\right) \\ D_1^S &= \frac{1}{M} \left(M - \frac{p_1 - p_2}{\Delta q}\right) \\ D_2^S &= \frac{1}{M} \left(\frac{p_1 - p_2}{\Delta q} - \frac{p_2}{q_2}\right) \end{aligned}$$

Thus, the total demand of product 1 and 2 are

$$\begin{aligned} D_1 &= (\beta + 1) - \left(\beta + \frac{1}{M}\right) \frac{p_1 - p_2}{\Delta q} \\ D_2 &= \left(\beta + \frac{1}{M}\right) \left[\frac{p_1 - p_2}{\Delta q} - \frac{p_2}{q_2}\right] \end{aligned}$$

(2) if $\frac{q_1}{p_1} > \frac{q_2}{p_2}$ (product 2 has higher unit cost of quality), then customers who have type $v \geq \frac{p_1}{q_1}$ purchase product 1; otherwise, they do not make any purchase. The total purchase from follower customers, D_1^F and D_2^F , and those from snob customers, D_1^S and D_2^S , are as follows, respectively.

$$\begin{aligned} D_1^F &= \beta \left(1 - \frac{p_1}{q_1}\right) \\ D_2^F &= 0 \\ D_1^S &= \frac{1}{M} \left(M - \frac{p_1}{q_1}\right) \\ D_2^S &= 0 \end{aligned}$$

Thus, the total demand of product 1 and 2 are

$$\begin{aligned} D_1 &= (\beta + 1) - \left(\beta + \frac{1}{M}\right) \frac{p_1}{q_1} \\ D_2 &= 0 \end{aligned}$$

Based on the demand functions above, we can derive the following response function of firm 1 and 2.

For firm 1, given any p_2 , (1) if he sets price $p_1 > \frac{p_2 q_1}{q_2}$, then the sales are $D_1 = (\beta + 1) - \left(\beta + \frac{1}{M}\right) \frac{p_1 - p_2}{\Delta q}$. In such case, the optimal price p_1^* is set to maximize his revenue $\pi_1 = p_1 D_1$: if $\frac{\beta + 1}{\beta + \frac{1}{M}} \Delta q > \left(\frac{2q_1}{q_2} - 1\right) p_2$, $p_1^* = \frac{1}{2} \left[p_2 + \frac{\beta + 1}{\beta + \frac{1}{M}} \Delta q \right]$; Otherwise, $p_1^* = \frac{p_2 q_1}{q_2}$. (2) If he sets price $p_1 \leq \frac{p_2 q_1}{q_2}$, then the sales is $D_1 = (\beta + 1) - \left(\beta + \frac{1}{M}\right) \frac{p_1}{q_1}$. In such cases, if $\frac{\beta + 1}{\beta + \frac{1}{M}} \leq \frac{2p_2}{q_2}$, $p_1^* = \frac{q_1}{2} \frac{\beta + 1}{\beta + \frac{1}{M}}$; Otherwise, $p_1^* = \frac{p_2 q_1}{q_2}$.

For firm 2, given any p_1 , (1) if he sets price $p_2 < \frac{p_1 q_2}{q_1}$, then the sales are $D_2 = \left(\beta + \frac{1}{M}\right) \left[\frac{p_1 - p_2}{\Delta q} - \frac{p_2}{q_2} \right]$. In such case, $p_2^* = \frac{p_1 q_2}{2q_1}$. (2) If he sets price $p_2 \geq \frac{p_1 q_2}{q_1}$, then his revenue is reduced to 0.

The best responses of firm 2 is to always choose to set price as $p_2^* = \frac{p_1 q_2}{2q_1}$. Firm 1, knowing the best response of firm 2, then set $p_1^* = \frac{\frac{1}{2} \frac{\beta + 1}{\beta + \frac{1}{M}} \Delta q}{1 - \frac{q_2}{4q_1}}$ (and firm 2 sets $p_2 = \frac{\beta + 1}{\beta + \frac{1}{M}} \frac{q_2}{4q_1 - q_2} \Delta q$).

Proof of Theorem III.8. In the first part of the analysis, we first start from the part where the demand expectation affects the customer purchases and optimal prices. Then, in the second part of the analysis, we analyze customers' expectations formed by the prices set by the firms.

Part 1: In the analysis, we also adopt the rational expectation framework. There are four cases in such a setting. (1) Snobs purchase product 1 (from firm 1) while followers purchase product 2 (from firm 2). (2) Both snobs and followers purchase product 1. (3) Both snobs and followers purchase product 2. (4) Snobs purchase product 2 while followers purchase product 1. Note that there does not exist cases where only snobs (or followers) makes a purchase (e.g., snobs purchase product 1 while followers make no purchase), as the no-sales firm can always attract the followers (or snobs) and be better off.

In case (1), $D_1^e = 1$ and $D_2^e = \beta$. Given p_2 , the price for firm 1 to get snob customer is $p_1 \leq v_S q_1 - k$ and $p_1 \leq v_S (q_1 - q_2) - k(1 - \beta) + p_2$ such that $v_S q_1 - k - p_1 \geq 0$ and

$$v_S q_1 - k - p_1 \geq v_S q_2 - k\beta - p_2.$$

$$p_1 = \begin{cases} v_S q_1 - k & \text{if } p_2 \geq v_S q_2 - k\beta \\ v_S(q_1 - q_2) - k(1 - \beta) + p_2 & \text{o.w.} \end{cases}$$

Given p_1 , the price for firm 2 to get follower customer is $p_2 \leq v_F q_2 + m\beta$ and $p_2 \leq p_1 - v_F(q_1 - q_2) - m(1 - \beta)$ such that $q_2 v_F + m\beta - p_2 \geq 0$ and $q_2 v_F + m\beta - p_2 \geq q_1 v_F + m - p_1$.

$$p_2 = \begin{cases} v_F q_2 + m\beta & \text{if } p_1 \geq v_F q_1 + m \\ p_1 - v_F(q_1 - q_2) - m(1 - \beta) & \text{o.w.} \end{cases}$$

By the best response functions of firm 1 and 2, we know that: (i) If $(v_S - v_F)\Delta q \geq (m + k)(1 - \beta)$, (a) if $(v_S - v_F)q_2 \geq \beta(m + k)$, $p_1 = v_S q_1 - k - (v_S - v_F)q_2 + \beta(k + m)$ and $p_2 = v_F q_2 + m\beta$; (b) Otherwise, $p_1 = v_S q_1 - k$ and $p_2 = v_F q_2 + m\beta$. (ii) If $(v_S - v_F)\Delta q < (m + k)(1 - \beta)$, then there does not exist equilibrium in this case.

In case (2), $D_1^e = 1 + \beta$ and $D_2^e = 0$. Given p_2 , the price for firm 1 to get snob customers is $p_1 \leq v_S q_1 - k(1 + \beta)$ and $p_1 \leq v_S(q_1 - q_2) - k(1 + \beta) + p_2$ such that $v_S q_1 - k(1 + \beta) - p_1 \geq 0$ and $v_S q_1 - k(1 + \beta) - p_1 \geq v_S q_2 - p_2$. And the price to get follower customers is $p_1 \leq v_F q_1 + m(1 + \beta)$ and $p_1 \leq v_F(q_1 - q_2) + m(1 + \beta) + p_2$ such that $v_F q_1 + m(1 + \beta) - p_1 \geq 0$ and $v_F q_1 + m(1 + \beta) - p_1 \geq v_F q_2 - p_2$. For firm 2, for any given p_1 , he has the incentive to lower p_2 to win customer. Thus, in the equilibrium, $p_1 = \min\{v_S(q_1 - q_2) - k(1 + \beta), v_F(q_1 - q_2) + m(1 + \beta)\}$ and $p_2 = 0$.

In case (3), $D_1^e = 0$ and $D_2^e = 1 + \beta$. Given p_1 , the price for firm 2 to get snob customers is $p_2 \leq v_S q_2 - k(1 + \beta)$ and $p_2 \leq p_1 - v_S(q_1 - q_2) - k(1 + \beta)$ such that $v_S q_2 - k(1 + \beta) - p_2 \geq 0$ and $v_S q_2 - k(1 + \beta) - p_2 \geq v_S q_1 - p_1$. And the price to get follower customers is $p_2 \leq v_F q_2 + m(1 + \beta)$ and $p_2 \leq p_1 - v_F(q_1 - q_2) + m(1 + \beta)$ such that $v_F q_2 + m(1 + \beta) - p_2 \geq 0$ and $v_F q_2 + m(1 + \beta) - p_2 \geq v_F q_1 - p_1$. For firm 1, for any given p_2 , he has the incentive to lower p_1 to win customer, until a point such that $p_2 = 0$ (from the constraints shown above). And from this point, any further reduce of p_1 means that firm 1 can generate positive sales, which conflicts with the demand expectation of $D_1^e = 0$. In other words, there does not exist an equilibrium that meet the demand expectations in this case.

In case (4), $D_1^e = \beta$ and $D_2^e = 1$. Given p_1 , the price for firm 1 to get follower customer is $p_1 \leq v_F q_1 + m\beta$ and $p_1 \leq v_F(q_1 - q_2) - m(1 - \beta) + p_2$ such that $v_F q_1 + m\beta - p_1 \geq 0$ and $v_F q_1 + m\beta - p_1 \geq v_F q_2 + m - p_2$. Given p_1 , the price for firm 2 to get snob customer is $p_2 \leq v_S q_2 - k$ and $p_2 \leq p_1 - v_S(q_1 - q_2) - k(1 - \beta)$ such that

$v_S q_2 - k - p_2 \geq 0$ and $v_S q_2 - k - p_2 \geq q_1 v_S - k\beta - p_1$. Note that such case is not feasible as the above constraints requires that $p_1 \leq v_F(q_1 - q_2) - m(1 - \beta) + p_2 \leq p_1 - m(1 - \beta) - (v_S - v_F)(q_1 - q_2) - k(1 - \beta)$, which does not hold.

Part 2: From the above analysis, for rational customers, if they observe that p_1 and p_2 are set as the same in case (1), they will form expectations of $D_1^e = 1$ and $D_2^e = \beta$; If they observe that p_1 and p_2 are set as the same in case (2), they will form expectations of $D_1^e = 0$ and $D_2^e = 1 + \beta$. In other words, firms pricing can fully determined customers' expectations and the corresponding purchases.

Next, consider firm 1's pricing decision. (i) If $(v_S - v_F)\Delta q \geq (m + k)(1 - \beta)$, (a) if $\min\{v_S q_1 - k - (v_S - v_F)q_2 - \beta(k + m), v_S q_1 - k\} \geq (1 + \beta) \min\{v_S(q_1 - q_2) - k(1 + \beta), v_F(q_1 - q_2) + m(1 + \beta)\}$, then the firm 1 prefers the equilibrium shown in case (1), and firm 2 also knows such preference. Thus, the equilibrium will be in the form shown in case (1). (b) Otherwise, firm 2 prefers the equilibrium shown in case (2), and firm 2 also knows such preference. Thus, the equilibrium will be in the form shown in case (2). (ii) If $(v_S - v_F)\Delta q < (m + k)(1 - \beta)$, then the only equilibrium is as the form shown in case (2).

Firms' Best Response Function. Based on the demand functions shown in Section 3.5.2, we can implicitly derive the following best response functions of firm 1 and 2. We organize our discussion in two scenarios, depending on the relation of quality and expected demand (as listed in the first column of conditions in Table 3.1).

(1) Scenario 1, $D_1 q_2 - D_2 q_1 \geq 0$.

For firm 1, for any given p_2 , he picks p_1 that maximizes profit π_1 .

$$\pi_1 = \left\{ \begin{array}{l} \max_{p_1} p_1 [(\beta + 1) - (\beta + \frac{1}{M})\frac{p_1}{q_1} + (\beta m - \frac{k}{M})\frac{D_1^e}{q_1}] \\ \text{if } p_1 \leq \frac{p_2 q_1 - k(D_1^e q_2 - D_2^e q_1)}{q_2}, \\ \max_{p_1} p_1 [(\beta + 1) - \beta \frac{-m D_1^e + p_1}{q_1} - \frac{1}{M} \frac{k(D_1^e - D_2^e) + p_1 - p_2}{\Delta q}] \\ \text{if } \frac{p_2 q_1 - k(D_1^e q_2 - D_2^e q_1)}{q_2} \leq p_1 \leq \frac{p_2 q_1 + m(D_1^e q_2 - D_2^e q_1)}{q_2}, \\ \max_{p_1} p_1 [(\beta + 1) - \beta \frac{-m(D_1^e - D_2^e) + p_1 - p_2}{\Delta q} - \frac{1}{M} \frac{k(D_1^e - D_2^e) + p_1 - p_2}{\Delta q}] \\ \text{if } p_1 \geq \frac{p_2 q_1 + m(D_1^e q_2 - D_2^e q_1)}{q_2}. \end{array} \right.$$

For firm 2, for any given p_1 , he picks the p_2 that maximizes profit π_2 .

$$\pi_2 = \begin{cases} 0, & \text{if } p_2 \geq \frac{p_1 q_2 + k(D_1^e q_2 - D_2^e q_1)}{q_1}, \\ \max_{p_2} p_2 \left[\frac{1}{M} \left(\frac{k(D_1^e - D_2^e) + p_1 - p_2}{\Delta q} - \frac{k D_2^e + p_2}{q_2} \right) \right], & \text{if } \frac{-m(D_1^e q_2 - D_2^e q_1) + p_1 q_2}{q_1} \leq p_2 \leq \frac{k(D_1^e q_2 - D_2^e q_1) + p_1 q_2}{q_1}, \\ \max_{p_2} p_2 \left[\beta \left(\frac{-m(D_1^e - D_2^e) + p_1 - p_2}{\Delta q} - \frac{-m D_1^e + p_1}{q_1} \right) + \frac{1}{M} \left(\frac{k(D_1^e - D_2^e) + p_1 - p_2}{\Delta q} - \frac{k D_2^e + p_2}{q_2} \right) \right], \\ \text{if } p_2 \leq \frac{-m(D_1^e q_2 - D_2^e q_1) + p_1 q_2}{q_1}. \end{cases}$$

(2) Scenario 2, $D_1^e q_2 - D_2^e q_1 < 0$.

For firm 1, for any given p_2 , he picks p_1 that maximizes profit π_1 .

$$\pi_1 = \begin{cases} \max_{p_1} p_1 \left[(\beta + 1) - \left(\beta + \frac{1}{M} \right) \frac{p_1}{q_1} + \left(\beta m - \frac{k}{M} \right) \frac{D_1^e}{q_1} \right] \\ \text{if } p_1 \leq \frac{p_2 q_1 + m(D_1^e q_2 - D_2^e q_1)}{q_2}, \\ \max_{p_1} p_1 \left[\beta + 1 - \beta \frac{-m(D_1^e - D_2^e) + p_1 - p_2}{\Delta q} - \frac{1}{M} \frac{k D_1^e + p_1}{q_1} \right] \\ \text{if } \frac{p_2 q_1 + m(D_1^e q_2 - D_2^e q_1)}{q_2} \leq p_1 \leq \frac{p_2 q_1 - k(D_1^e q_2 - D_2^e q_1)}{q_2}, \\ \max_{p_1} p_1 \left[(\beta + 1) - \beta \frac{-m(D_1^e - D_2^e) + p_1 - p_2}{\Delta q} - \frac{1}{M} \frac{k(D_1^e - D_2^e) + p_1 - p_2}{\Delta q} \right] \\ \text{if } p_1 \geq \frac{p_2 q_1 - k(D_1^e q_2 - D_2^e q_1)}{q_2}. \end{cases}$$

For firm 2, for any given p_1 , he picks the p_2 that maximizes profit π_2 .

$$\pi_2 = \begin{cases} 0, & \text{if } p_2 \geq \frac{-m(D_1^e q_2 - D_2^e q_1) + p_1 q_2}{q_1}, \\ \max_{p_2} p_2 \left[\beta \left(\frac{-m(D_1^e - D_2^e) + p_1 - p_2}{\Delta q} - \frac{-m D_1^e + p_1}{q_1} \right) \right], & \text{if } \frac{k(D_1^e q_2 - D_2^e q_1) + p_1 q_2}{q_1} \leq p_2 \leq \frac{-m(D_1^e q_2 - D_2^e q_1) + p_1 q_2}{q_1}, \\ \max_{p_2} p_2 \left[\beta \left(\frac{-m(D_1^e - D_2^e) + p_1 - p_2}{\Delta q} - \frac{-m D_1^e + p_1}{q_1} \right) + \frac{1}{M} \left(\frac{k(D_1^e - D_2^e) + p_1 - p_2}{\Delta q} - \frac{k D_2^e + p_2}{q_2} \right) \right], \\ \text{if } p_2 \leq \frac{k(D_1^e q_2 - D_2^e q_1) + p_1 q_2}{q_1}. \end{cases}$$

Proofs of Section 3.6

Proof of Lemma III.12.

From Assumption III.3, we know that there are three possible scenarios: (1) Only snob customers purchase the bundle. (2) Both snob and follower customers purchase the bundle. (3) Only follower customers purchase the bundle.

(1) In the first scenario, $D_1^e = 1$ and $D_2^e = 1$. It is easy to see that the optimal bundle price is $p_B = v_S q_1 + v_S q_2 - 2k$. The revenue is $R_B = p_B = p_B = v_S q_1 + v_S q_2 - 2k$.

(2) In the second scenario, $D_1^e = 1 + \beta$ and $D_2^e = 1 + \beta$. The optimal bundle price is $p_B = v_F q_1 + v_F q_2 + 2m(1 + \beta)$. Revenue is $R_B = (1 + \beta)[v_F q_1 + v_F q_2 + 2m(1 + \beta)]$.

(3) In the third scenario, $D_1^e = \beta$ and $D_2^e = \beta$. The optimal bundle price is $p_B = v_F q_1 + v_F q_2 + 2\beta m$. Note that such a case appears only if $v_F q_1 + v_F q_2 + 2\beta m \geq v_S q_1 + v_S q_2 - 2k$. Revenue is $R_B = \beta(v_F q_1 + v_F q_2 + 2\beta m)$. Note that the revenue is dominated by that in the second scenario.

Thus, the revenue from pure bundling is $\max\{v_S q_1 + v_S q_2 - 2k, (1 + \beta)[v_F q_1 + v_F q_2 + 2m(1 + \beta)]\}$.

Proof of Lemma III.14

We analyze this case by checking the 2 possible strategies of product bundling: (1) Offering only the bundle and product 2. (2) Offering only the bundle and product 1.

Regarding the first strategy, there are two scenarios of the realized sales. We analyze them one by one as follows.

Scenario 1: Snob customers purchase bundle and follower customers purchase product 2. $D_1^e = 1$ and $D_2^e = 1 + \beta$.

The firm chooses price p_B and p_2 to maximize profit, with the above set of constraints.

$$\begin{aligned}
 \max \quad & p_B + p_2 \beta \\
 \text{s.t.} \quad & v_S q_1 - k \geq p_B - p_2 \\
 & v_S q_1 + v_S q_2 - k - p_B \geq 0 \\
 & v_F q_1 + m \leq p_B - p_2 \\
 & v_F q_2 + m(1 + \beta) - p_2 \geq 0
 \end{aligned}$$

where the constraints represents the conditions for snobs to purchase the bundle and only snobs want to purchase the bundle. First, the utility from bundle should be larger than that of product 2 for snobs $v_S q_1 + v_S q_2 - k - p_B \geq v_S q_2 - p_2 \iff v_S q_1 - k \geq p_B - p_2$.

Second, the utility of bundle is positive. Third, for followers, bundle utility is smaller than separate product 2, $v_F(q_1 + q_2) + m + m(1 + \beta) - p_B < v_F q_2 + m(1 + \beta) - p_2 \iff v_F q_1 + m < p_B - p_2$. Last, product 2 has positive utility for follower customer.

Note that the above optimization is feasible if $v_S q_1 - k \geq v_F q_1 + m$, which automatically holds with Assumption III.10.

Solving the above optimization, the optimal prices are

$$\begin{aligned} p_2^* &= \min\{v_F q_2 + m(1 + \beta), v_S(q_1 + q_2) - v_F q_1 - k - m\} \\ p_B^* &= v_S(q_1 + q_2) - k \end{aligned}$$

One observation is that, depending on the paramters, the snob customers may have positive utility, while follower customers always have zero utility.

The total revenue is then

$$R = v_S(q_1 + q_2) - k + \beta \min\{v_F q_2 + m(1 + \beta), v_S(q_1 + q_2) - v_F q_1 - k - m\}$$

Scenario 2: Snob customers only purchase 2 and follower customers only purchase the bundle. $D_1^e = \beta$ and $D_2^e = 1 + \beta$.

The following conditions must hold in this case. First, for snobs customers, the utility of bundle is smaller than 2. $v_S q_1 + v_S q_2 - k\beta - p_B \leq v_S q_2 - p_2 \iff v_S q_1 - k\beta \leq p_B - p_2$. And, for followers, the utility of product 2 is smaller than that of the bundle. $v_F q_1 + v_F q_2 + \beta m + m(1 + \beta) - p_B > v_F q_2 + m(1 + \beta) - p_2 \iff v_F q_1 + \beta m > p_B - p_2$. By Assumption III.10, it is easy to see that there is no feasible solutions in this case.

Regarding the second strategy, by similar logic, we can derive that the total revenue is

$$R = v_S(q_1 + q_2) - k(1 + \beta) + \beta \min\{v_F q_1 + m(1 + \beta), v_S(q_1 + q_2) - v_F q_2 - k - m\}$$

Thus, the optimal selling strategy is to offer bundle and product 1 if $\min\{v_F q_1 + m(1 + \beta), v_S(q_1 + q_2) - v_F q_2 - k - m\} - \min\{v_F q_2 + m(1 + \beta), v_S(q_1 + q_2) - v_F q_1 - k - m\} \leq k$. Note that this condition can be easily simplified as $v_F(q_1 - q_2) \leq k$; Otherwise, the optimal strategy is to offer bundle and product 2.

Proof of Lemma III.15

For the cases of $k = m = 0$, from propositions III.11, III.12, III.13, and III.14, we

know that the revenue of each selling strategy is as below

$$\begin{aligned}
R_{Pure\ Bundling} &= \max\{v_S(q_1 + q_2), (1 + \beta)v_F(q_1 + q_2)\} \\
R_{Pure\ Component} &= \max\{v_Sq_1, (1 + \beta)v_Fq_1\} + \max\{v_Sq_2, (1 + \beta)v_Fq_2\} \\
R_{Partial\ Mixed\ Bundle} &= v_S(q_1 + q_2) + \beta \min\{v_Fq_1, v_S(q_1 + q_2) - v_Fq_2\}
\end{aligned}$$

For the cases of $v_S < (1 + \beta)v_F$, the revenues from pure bundling, pure component, and partial mixed bundling are $(1 + \beta)v_F(q_1 + q_2)$, $v_Sq_1 + (1 + \beta)v_Fq_2$, and $v_S(q_1 + q_2) + \beta v_Fq_1$, respectively. Note that the partial mixed bundling and the pure bundling cases always dominate the pure component strategy. If $(v_S - v_F)(q_1 + q_2) \geq \beta v_Fq_2$, then the mixed bundling dominates the pure bundling strategy; Otherwise, the pure bundling strategy is optimal.

For the cases of $v_S \geq (1 + \beta)v_F$, the revenues from pure bundling, pure component, and partial mixed bundling are $v_S(q_1 + q_2)$, $v_Sq_1 + (1 + \beta)v_Fq_2$, and $v_S(q_1 + q_2) + \beta v_Fq_1$, respectively. It is easy to see that the partial mixed bundling always dominates other strategies.

Proof of Lemma III.16

We identify the cases where partial mixed bundling outperform the other strategies by analyzing the subcases one by one.

Scenario (1): $k \geq v_F(q_1 - q_2)$. Remember that in the cases with $k \geq v_F(q_1 - q_2)$, offering product 2 with bundle is better than offering product 1 with bundle. The revenue of each strategy in as below.

$$\begin{aligned}
R_{Pure\ Bundling} &= \max\{v_S(q_1 + q_2) - k, (1 + \beta)[v_F(q_1 + q_2) + 2m(1 + \beta)]\} \\
R_{Pure\ Component} &= \max\{v_Sq_1 - k, (1 + \beta)[v_Fq_1 + m(1 + \beta)]\} \\
&\quad + \max\{v_Sq_2, (1 + \beta)[v_Fq_2 + m(1 + \beta)]\} \\
R_{Partial\ Mixed\ Bundle} &= v_S(q_1 + q_2) - k \\
&\quad + \beta \min\{v_Fq_2 + m(1 + \beta), v_S(q_1 + q_2) - v_Fq_1 - k - m\}
\end{aligned}$$

Note that if **(a)** $v_Sq_1 - k \geq (1 + \beta)[v_Fq_1 + m(1 + \beta)]$ and $v_Sq_2 \geq (1 + \beta)[v_Fq_2 + m(1 + \beta)]$ or **(b)** $v_Sq_1 - k < (1 + \beta)[v_Fq_1 + m(1 + \beta)]$ and $v_Sq_2 < (1 + \beta)[v_Fq_2 + m(1 + \beta)]$, the revenue of pure component is equal to that of pure bundling. The only case where they are not equal is $v_Sq_1 - k \geq (1 + \beta)[v_Fq_1 + m(1 + \beta)]$ and $v_Sq_2 < (1 + \beta)[v_Fq_2 + m(1 + \beta)]$. In such a case, $R_{Pure\ Component} = v_Sq_1 - k + (1 + \beta)[v_Fq_2 + m(1 + \beta)]$. In the following

analysis, with a bit abuse of notation, we only refer the pure component to this case.

(1) If $v_S(q_1 + q_2) - k < (1 + \beta)[v_F(q_1 + q_2) + 2m(1 + \beta)]$ and $(v_S - v_F)(q_1 + q_2) < (2 + \beta)m + k$, pure bundling revenue is $R_{Pure\ Bundle} = (1 + \beta)[v_F(q_1 + q_2) + 2m(1 + \beta)]$ and the partial mixed bundling is $R_{Partial\ Mixed\ Bundle} = v_S(q_1 + q_2) - k + \beta[v_S q_2 + (v_S - v_F)q_1 - k - m]$. The partial mixed bundling dominates pure bundling if $(1 + \beta)(v_S - v_F)(q_1 + q_2) - \beta v_F q_1 \geq k(1 + \beta) + m\beta + 2m(1 + \beta)^2$, which contradicts the second condition in (1). Thus, the partial mixed bundling is always dominated by pure bundling in this case.

(2) If $v_S(q_1 + q_2) - k < (1 + \beta)[v_F(q_1 + q_2) + 2m(1 + \beta)]$ and $(v_S - v_F)(q_1 + q_2) \geq (2 + \beta)m + k$, pure bundling revenue is $R_{Pure\ Bundle} = (1 + \beta)[v_F(q_1 + q_2) + 2m(1 + \beta)]$ and the partial mixed bundling revenue is $R_{Partial\ Mixed\ Bundle} = v_S(q_1 + q_2) - k + \beta[v_F q_2 + m(1 + \beta)]$. The partial mixed bundling dominates pure bundling if $(v_S - v_F)(q_1 + q_2) - \beta v_F q_1 \geq (2 + \beta)(1 + \beta)m + k$. And it dominates the pure component if $(v_S - v_F)q_2 \geq m(1 + \beta)$, which always hold by Assumption III.10.

(3) If $v_S(q_1 + q_2) - k \geq (1 + \beta)[v_F(q_1 + q_2) + 2m(1 + \beta)]$ and $(v_S - v_F)(q_1 + q_2) < (2 + \beta)m + k$, these two conditions conflict each other. So there does not exist such cases.

(4) If $v_S(q_1 + q_2) - k \geq (1 + \beta)[v_F(q_1 + q_2) + 2m(1 + \beta)]$ and $(v_S - v_F)(q_1 + q_2) \geq (2 + \beta)m + k$, pure bundling revenue is $R_{Pure\ Bundle} = v_S(q_1 + q_2) - k$ and the partial mixed bundling revenue is $R_{Partial\ Mixed\ Bundle} = v_S(q_1 + q_2) - k + \beta[v_F q_2 + m(1 + \beta)]$. The partial mixed bundling dominates the pure bundling if $v_F q_2 + m(1 + \beta) \geq k$. And it dominated the pure component if $(v_S - v_F)q_2 \geq m(1 + \beta)$, which always hold by Assumption III.10.

Scenario (2): $k < v_F(q_1 - q_2)$, where offering product 1 with bundle is better than offering product 2 with bundle. The revenue of each strategy in as below.

$$\begin{aligned}
R_{Pure\ Bundling} &= \max\{v_S(q_1 + q_2) - k, (1 + \beta)[v_F(q_1 + q_2) + 2m(1 + \beta)]\} \\
R_{Pure\ Component} &= \max\{v_S q_1 - k, (1 + \beta)[v_F q_1 + m(1 + \beta)]\} \\
&\quad + \max\{v_S q_2, (1 + \beta)[v_F q_2 + m(1 + \beta)]\} \\
R_{Partial\ Mixed\ Bundle} &= v_S(q_1 + q_2) - k(1 + \beta) \\
&\quad + \beta \min\{v_F q_1 + m(1 + \beta), v_S(q_1 + q_2) - v_F q_2 - k - m\}
\end{aligned}$$

Again if (a) $v_S q_1 - k \geq (1 + \beta)[v_F q_1 + m(1 + \beta)]$ and $v_S q_2 \geq (1 + \beta)[v_F q_2 + m(1 + \beta)]$ or (b) $v_S q_1 - k < (1 + \beta)[v_F q_1 + m(1 + \beta)]$ and $v_S q_2 < (1 + \beta)[v_F q_2 + m(1 + \beta)]$, the

revenue of pure component is equal to that of pure bundling. The only case where they are not equal is $v_S q_1 - k \geq (1 + \beta)[v_F q_1 + m(1 + \beta)]$ and $v_S q_2 < (1 + \beta)[v_F q_2 + m(1 + \beta)]$. In such a case, $R_{Pure\ Component} = v_S q_1 - k + (1 + \beta)[v_F q_2 + m(1 + \beta)]$. In the following analysis, with a bit abuse of notation, we only refer the pure component to this case.

(1) If $v_S(q_1 + q_2) - k < (1 + \beta)[v_F(q_1 + q_2) + 2m(1 + \beta)]$ and $(v_S - v_F)(q_1 + q_2) < (2 + \beta)m + k$, pure bundling revenue is $R_{Pure\ Bundle} = (1 + \beta)[v_F(q_1 + q_2) + 2m(1 + \beta)]$ and the partial mixed bundling is $R_{Partial\ Mixed\ Bundle} = v_S(q_1 + q_2) - k(1 + \beta) + \beta[v_S(q_1 + q_2) - v_F q_2 - k - m]$. The partial mixed bundling dominates pure bundling if $(1 + \beta)(v_S - v_F)(q_1 + q_2) - \beta v_F q_2 \geq k(1 + 2\beta) + m\beta + 2m(1 + \beta)^2$, which contradicts with the second condition in (1) (of Scenario (2)). Thus, partial mixed bundling is always dominated by pure bundling in this case.

(2) If $v_S(q_1 + q_2) - k < (1 + \beta)[v_F(q_1 + q_2) + 2m(1 + \beta)]$ and $(v_S - v_F)(q_1 + q_2) \geq (2 + \beta)m + k$, pure bundling revenue is $R_{Pure\ Bundle} = (1 + \beta)[v_F(q_1 + q_2) + 2m(1 + \beta)]$ and the partial mixed bundling revenue is $R_{Partial\ Mixed\ Bundle} = v_S(q_1 + q_2) - k(1 + \beta) + \beta[v_F q_1 + m(1 + \beta)]$. The partial mixed bundling dominates pure bundling if $(v_S - v_F)(q_1 + q_2) - \beta v_F q_2 \geq (2 + \beta)(1 + \beta)m + k(1 + \beta)$. And it dominates the pure component if $(v_S - v_F)q_2 + \beta v_F(q_1 - q_2) \geq m(1 + \beta) + k\beta$, which always hold by Assumption III.10.

(3) If $v_S(q_1 + q_2) - k \geq (1 + \beta)[v_F(q_1 + q_2) + 2m(1 + \beta)]$ and $(v_S - v_F)(q_1 + q_2) < (2 + \beta)m + k$, these two conditions conflict each other. So there does not exist such cases.

(4) If $v_S(q_1 + q_2) - k \geq (1 + \beta)[v_F(q_1 + q_2) + 2m(1 + \beta)]$ and $(v_S - v_F)(q_1 + q_2) \geq (2 + \beta)m + k$, pure bundling revenue is $R_{Pure\ Bundle} = v_S(q_1 + q_2) - k$ and the partial mixed bundling revenue is $R_{Partial\ Mixed\ Bundle} = v_S(q_1 + q_2) - k(1 + \beta) + \beta[v_F q_1 + m(1 + \beta)]$. The partial mixed bundling dominates the pure bundling if $v_F q_1 + m(1 + \beta) \geq k$. And it dominated the pure component if $(v_S - v_F)q_2 + \beta v_F(q_1 - q_2) \geq m(1 + \beta) + k\beta$, which always hold by Assumption III.10.

BIBLIOGRAPHY

- [1] Acimovic, J. and S.C. Graves. 2015. Making better fulfillment decisions on the fly in an online retail environment. *Manufacturing & Service Operations Management* **17** (1) 34-51.
- [2] Adams, W. J., and Yellen, J. L. 1976. Commodity bundling and the burden of monopoly. *The Quarterly Journal of Economics*, 475-498.
- [3] Ahn, H.S., S. Jasin, P. Kaminsky, Y. Wang. 2018. Analysis of Deterministic Control and Its Improvements for an Inventory Problem with Multiproduct Batch Differentiation. *Operations Research*. **66**(1) 58-76.
- [4] Amaldoss, W., and Jain, S. 2005. Conspicuous consumption and sophisticated thinking. *Management Science*, **51**(10), 1449-1466.
- [5] Amaldoss, W., and Jain, S. 2005. Pricing of conspicuous goods: A competitive analysis of social effects. *Journal of Marketing Research*, **42**(1), 30-42.
- [6] Amaldoss, W., and Jain, S. 2008. Research note—Trading up: A strategic analysis of reference group effects. *Marketing Science*, **27**(5), 932-942.
- [7] Amaldoss, Wilfred, and Sanjay Jain. 2015. Branding conspicuous goods: An analysis of the effects of social influence and competition. *Management Science*. **61** (9), 2064-2079.
- [8] Armstrong, M., and Vickers, J. 2010. Competitive non-linear pricing and bundling. *The Review of Economic Studies*, **77**(1), 30-60.
- [9] ASAP Expedited Logistics. 2018. www.asapexpediting.net.
- [10] Balachander, S., and Stock, A. 2009. Limited edition products: When and when not to offer them. *Marketing Science*. **28** (2), 336-355.
- [11] Bain and Company. 2019. Luxury goods worldwide market study, Fall-Winter 2018.
- [12] Berghaus, B., Müller-Stewens, G., and Reinecke, S. 2014. The management of luxury: A practitioner's handbook. Kogan Page Publishers.
- [13] Bertsimas, D., I. Ch. Paschalidis. 2001. Probabilistic Service Level Guarantees in Make-to-Stock Manufacturing Systems. *Operations Research*. **49**(1) 119-133.
- [14] Bijvank, M. 2014. Periodic Review Inventory System with a Service Level Criterion. *Journal of the Operational Research Society*. **65** 1853-1863.
- [15] Bijvank, M., I.F.A. Vis. 2011. Lost-sales Inventory Theory: A Review. *European Journal of Operations Research*. **215** 1-13.
- [16] Bijvank, M., I.F.A. Vis. 2012. Lost-sales Inventory System with a Service Level Criterion. *European Journal of Operations Research*. **220** 610-618.

- [17] Bitran, G. R., T.-Y. Leong. 1992. Deterministic Approximations to Co-Production Problems with Service Constraints and Random Yields. *Management Science*. **38**(5) 724-742.
- [18] Bitran, G.R., D. Sarkar. 1988. On Upper Bounds of Sequential Stochastic Production Planning Problems. *European Journal of Operational Research*. **34**(2) 191-207.
- [19] Bitran, G. R., H. H. Yanasse. 1984. Deterministic Approximations to Stochastic Production Problems. *Operations Research*. **32**(5) 999-1018.
- [20] Bloomberg, June 2015. <https://www.bloomberg.com/news/articles/2015-06-10/how-the-legendary-birkin-bag-remains-dominant>.
- [21] Boyaci, T., G. Gallego. 2001. Serial Production/Distribution Systems Under Service Constraints. *Manufacturing and Service Operations Management*. **3**(1) 43-50.
- [22] Bu, Jinzhi, Xiting Gong, Dacheng Yao. 2017. Constant-order Policies for Lost-sales Inventory Models with Random Supply Functions: Asymptotics and Heuristic. Available at SSRN: <https://ssrn.com/abstract=3063730>.
- [23] Burns, L.D., R.W. Hall, D.E. Blumenfeld, C.F. Daganzo. 1985. Distribution strategies that minimize transportation and inventory costs. *Operations Research* **33** (3) 469-490.
- [24] Cachon, G. P., and Swinney, R. 2009. Purchasing, pricing, and quick response in the presence of strategic consumers. *Management Science*, **55**(3), 497-511.
- [25] Caggiano, K.E., J.A Muckstadt, J.A. Rappold. 2006. Integrated real-time capacity and inventory allocation for reparable service parts in a two-echelon supply system. *Manufacturing & Service Operations Management* **8** (3) 292-319.
- [26] Campbell, J.F. 1990. Designing logistics systems by analyzing transportation, inventory and terminal cost tradeoffs. *Journal of Business Logistics* **11** (1) 159-179.
- [27] Candogan, O., Bimpikis, K., and Ozdaglar, A. 2012. Optimal pricing in networks with externalities. *Operations Research*, **60** (4): 883-905.
- [28] Cetinkaya, S. 2005. Coordination of inventory and shipment consolidation decisions: A review of premises, models, and justification. In *Applications of supply chain management and e-commerce research* (pp. 3-51). Springer US.
- [29] Cetinkaya, S., E. Tekin, C.Y. Lee. 2008. A stochastic model for integrated inventory replenishment and outbound shipment release decisions. *IIE Transactions* **40** (8) 324-340.
- [30] Chen, F. Y., D. Krass. 2001. Inventory Models with Minimal Service Level Constraints. *European Journal of Operational Research*. **134**(1) 120-140.
- [31] Corneo, G., and Jeanne, O. 1997. Conspicuous consumption, snobbism and conformism. *Journal of Public Economics*, **66**(1), 55-71.
- [32] Daganzo, C.F. 1988. Shipment composition enhancement at a consolidation center. *Transportation Research – Part B* **22** (2) 103-124.

- [33] Dittmar, Helga. 1994. Material possessions as stereotypes: Material images of different socio-economic groups. *Journal of Economic Psychology* **15** (4), 561-585.
- [34] Economist, September 2018. <https://www.economist.com/business/2018/09/27/makers-of-very-expensive-cars-want-to-be-luxury-goods-firms?fsrc=scn/fb/te/bl/ed/makersofvery-expensivecarswanttobeluxurygoodsfirmsjoiningthehighrevvers>.
- [35] Expedited Logistics and Freight Services. 2018. www.elfsfreight.com/services-domestic.php
- [36] FedEx Standard List Rates. 2018. images.fedex.com/us/services/pdf/FedEx_Standard_ListRates_2018.pdf
- [37] Forbes, January 2018. <https://www.forbes.com/sites/gracelwilliams/2016/01/16/sell-stocks-buy-a-birkin-not-so-fast-experts-say>.
- [38] Ghosh, B., and Balachander, S. 2007. Research note—competitive bundling and counterbundling with generalist and specialist firms. *Management Science*, **53**(1), 159-168.
- [39] Gil, R., E. Korkmaz, and O. Sahin. 2014. Optimal Pricing of Access and Secondary Goods with Repeat Purchases: Evidence from Online Grocery Shopping and Delivery Fees. NET Institute Working Paper No. 14-10.
- [40] Goldberg, D.A., D. A. Katz-Rogozhnikov, Y. Lu, M. Sharma, M. S. Squillante. 2016. Asymptotic Optimality of Constant-Order Policies for Lost Sales Inventory Models with Large Lead Times. *Mathematics of Operations Research*. **41**(3) 898-913.
- [41] Gupta, Y.P., P.K. Bagchi. 1987. Inbound freight consolidation under just-in-time procurement: application of clearing models. *Journal of Business Logistics* **8** (2) 74-94.
- [42] Gurzki, H., and Woisetschläger, D. M. 2017. Mapping the luxury research landscape: A bibliometric citation analysis. *Journal of Business Research*, **77**, 147-166.
- [43] Hausmann, L., N.A. Herrmann, J. Krause, and T. Netzer. 2014. Same-day delivery: The next evolutionary step in parcel logistics. Retrieved from <http://www.mckinsey.com/industries/travel-transport-and-logistics/our-insights/same-day-delivery-the-next-evolutionary-step-in-parcel-logistics>.
- [44] Higginson, J.K., and J.H. Bookbinder. 1994. Policy recommendations for a shipment consolidation program. *Journal of Business Logistics*. **15** (1) 87-112.
- [45] Hoadley, B. and D.P. Heyman. 1977. A two—echelon inventory model with purchases, dispositions, shipments, returns and transshipments. *Naval Research Logistics Quarterly* **24** (1) 1-19.
- [46] Huggins, E.L., T.L. Olsen. 2003. Supply chain management with guaranteed delivery. *Management Science* **40** (9) 1154-1167.
- [47] Huh, W.T., G. Janakiraman, J.A. Muckstadt, P. Rusmevichientong. 2009. Asymptotic Optimality of Order-Up-To Policies in Lost Sales Inventory Systems. *Management Science*. **55**(3) 404-420.

- [48] Janakiraman, G., S. Seshadri, G. Shanthikumar. 2007. A Comparison of the Optimal Costs of Two Canonical Inventory Systems. *Operations Research*. **55**(5) 866-875.
- [49] Jasin, S., and A. Sinha. 2015. An LP-based correlated rounding scheme for multi-item ecommerce order fulfillment. *Operations Research* **63** (6) 1336-1351.
- [50] Jiang, Y., C. Shi, S. Shen. 2019. Service-Level Constrained Inventory Systems. *Production and Operations Management*. Forthcoming.
- [51] Jing, B. 2007. Network externalities and market segmentation in a monopoly. *Economics Letters*. **95** (1): 7-13.
- [52] Katz, M. L., and Shapiro, C. 1985. Network externalities, competition, and compatibility. *American Economic Review*. **75** (3): 424-440.
- [53] Kapferer, Jean-Noël. 2016. Kapferer on luxury: How luxury brands can grow yet remain rare. Kogan Page Publishers.
- [54] Laffont, J. J., Rey, P., and Tirole, J. 1998. Network competition: I. Overview and nondiscriminatory pricing. *The RAND Journal of Economics*, 1-37.
- [55] Lei Y., S. Jasin and A. Sinha. 2016. Dynamic joint pricing and order fulfillment for e-commerce retailers. (March 23, 2016). Ross School of Business Paper No. 1310.
- [56] Leibenstein, H. 1950. Bandwagon, snob, and Veblen effects in the theory of consumers' demand. *The Quarterly Journal of Economics*. **64**(2), 183-207.
- [57] Lewis, M., V. Singh, and S. Fay. 2006. An empirical study of the impact of nonlinear shipping and handling fees on purchase incidence and expenditure decisions. *Marketing Science* **25** (1): 51-64.
- [58] Lovasz, L. 1993. *Combinatorial problems and exercises*. Vol. 361. American Mathematical Soc.
- [59] Ma, M., and Mallik, S. 2017. Bundling of vertically differentiated products in a supply chain. *Decision Sciences*, **48**(4), 625-656.
- [60] Matutes, C., and Regibeau, P. 1988. "Mix and match": product compatibility without network externalities. *The RAND Journal of Economics*, 221-234.
- [61] Momot, R., Belavina, E., and Girotra, K. 2019. The use and value of social network information in selective selling. *Management Science*, forthcoming.
- [62] Ng, C.K. Inbound supply chain optimization and process improvement. Diss. Massachusetts Institute of Technology, 2012.
- [63] Özer, Ö., H. Xiong. 2008. Stock positioning and performance estimation for distribution systems with service constraints. *IIE Transactions*. **40**(12) 1141-1157.
- [64] Prasad, A., Venkatesh, R., and Mahajan, V. 2010. Optimal bundling of technological products with network externality. *Management Science*. **56** (12): 2224-2236.
- [65] Qi, L., and K. Lee. Supply chain risk mitigations with expedited shipping. *Omega* **57** : 98-113.

- [66] Rao, R. S., and Schaefer, R. 2013. Conspicuous consumption and dynamic pricing. *Marketing Science*, **32(5)**, 786-804.
- [67] Reiman, M. I., Q. Wang. Asymptotically Optimal Inventory Control for Assemble-To-Order Systems With Identical Lead Times. *Operations Research*. **63** 3 716-732.
- [68] Shang, K. H., J. S. Song. 2006. A Closed-Form Approximation for Serial Inventory Systems and Its Application to System Design. *Manufacturing and Service Operations Management*. **8(4)** 394-406.
- [69] Schmalensee, R. 1984. Gaussian demand and commodity bundling. *Journal of Business*, S211-S230.
- [70] Shaked, A., and Sutton, J. 1982. Relaxing price competition through product differentiation. *The Review of Economic Studies*, 3-13.
- [71] Snyder, L. V., Z.-J. M. Shen. 2011. *Fundamentals of supply chain theory*. John Wiley & Sons.
- [72] Stock, A., and Balachander, S. 2005. The making of a “hot product”: A signaling explanation of marketers’ scarcity strategy. *Management Science*, **51(8)**, 1181-1192.
- [73] Stremersch, S., and Tellis, G. J. 2002. Strategic bundling of products and prices: A new synthesis for marketing. *Journal of Marketing*, **66(1)**, 55-72.
- [74] Su, X., and Zhang, F. 2008. Strategic customer behavior, commitment, and supply chain performance. *Management Science*, **54(10)**, 1759-1773.
- [75] Tarim, S. A., B. G. Kingsman. 2003. The Stochastic Dynamic Production/Inventory Lot-Sizing Problem with Service Level Constraints. *International Journal of Production Economics*. **88** 105-119.
- [76] Tarim, S. A., M. K. Dogru, U. Ozen, R. Rossi. 2011. An Efficient Computational Method for a Stochastic Dynamic Lot-Sizing Problem Under Service Level Constraints. *European Journal of Operations Research*. **215** 563-571.
- [77] Tereyağoğlu, N., and Veeraraghavan, S. 2012. Selling to conspicuous consumers: Pricing, production, and sourcing decisions. *Management Science*, **58(12)**, 2168-2189.
- [78] UPS. 2018. www.ups.com/us/en/shipping/zones-and-rates/48-contiguous-states.page
- [79] USPS. 2018. <https://pe.usps.com/text/dmm300/Notice123.htm>
- [80] Veblen, Thorstein. 1899. The theory of the leisure class: an economic study of institutions. New York: The Modern Library (1934).
- [81] Veeraraghavan, S., and Debo, L. 2009. Joining longer queues: Information externalities in queue choice. *Manufacturing & Service Operations Management*, **11(4)**, 543-562.
- [82] Veeraraghavan, S. K., and Debo, L. G. 2011. Herding in queues with waiting costs: Rationality and regret. *Manufacturing & Service Operations Management*, **13(3)**, 329-346.

- [83] Vigneron, F., and Johnson, L. W. 1999. A review and a conceptual framework of prestige-seeking consumer behavior. *Academy of Marketing Science Review*. **1 (1)**: 1-15.
- [84] Viswanathan, S. 2005. Competing across technology-differentiated channels: The impact of network externalities and switching costs. *Management Science*, **51(3)**, 483-496.
- [85] The Wall Street Journal. January 2007. www.wsj.com/articles/SB116836324469271556.
- [86] Wan, H., Q. Wang. Asymptotically-Optimal Component Allocation for Assemble-To-Order Production-Inventory Systems. *Operations Research Letters*. **43(3)** 304-310.
- [87] Wilson, R. 2016. Council of Supply Chain Management Professionals. *27th Annual State of Logistics Report*.
- [88] Worm, K. G. J. P., E. M. T. Hendrix, R. Haijema, J. G.A.J van der Vorst. 2014. An MILP Approximation for Ordering Perishable Products with Non-stationary Demand and Service Level Constraints. *International Journal of Production Economics*. **157** 133-146.
- [89] Xin, L., D.A. Goldberg. 2016. Optimality Gap of Constant-Order Policies Decays Exponentially in the Lead Time for Lost Sales Models. *Operations Research*. **64(6)** 1556-1565.
- [90] Xin, Linwei, Long He, Jagtej Bewli, John Bowman, Huijun Feng, Zhiwei (Tony) Qin. 2017. On the Performance of Tailored Base-Surge Policies: Theory and Application at Walmart.com. Available at SSRN: <https://ssrn.com/abstract=3090177>.
- [91] Xin, L. 2019. Understanding the Performance of Capped Base-Stock Policies in Lost-Sales Inventory Models. Available at SSRN: <https://ssrn.com/abstract=3357241>.
- [92] Xu, P.J., R. Allgor, and S.C. Graves. 2009. Benefits of reevaluating real-time order fulfillment decisions. *Manufacturing & Service Operations Management* **11 (2)** 340-355.
- [93] Zhou, S.X., and X. Chao. 2010. Newsvendor bounds and heuristics for serial supply chains with regular and expedited shipping. *Naval Research Logistics* **57 (1)** : 71-87.
- [94] Zhou, J. 2017. Competitive bundling. *Econometrica*, **85(1)**, 145-172.
- [95] Zipkin, P. 2008. Old and New Methods for Lost-Sales Inventory Systems. *Operations Research*. **56(5)** 1256-1263.
- [96] Zipkin, P. 2000. *Foundations of Inventory Management*. Vol. 2. New York, McGraw-Hill.