# Information and Decentralization in Inventory, Supply Chain, and Transportation Systems

by

Guillaume Roels

Submitted to the Sloan School of Management
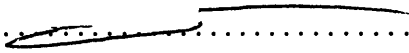in partial fulfillment of the requirements for the degree of
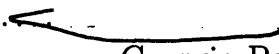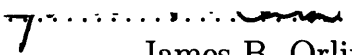
Doctor of Philosophy in Operations Research

at the
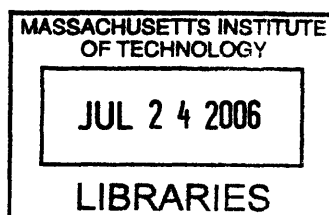
MASSACHUSETTS INSTITUTE OF TECHNOLOGY

June 2006

Author . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .
Sloan School of Management
May 18, 2006

Certified by. . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .
Georgia Perakis
J. Spencer Standish Associate Professor of Operations Research
Sloan School of Management
Thesis Supervisor

Accepted by . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .
James B. Orlin
Edward Pennell Brooks Professor of Operations Research
Codirector, Operations Research Center

# Information and Decentralization in Inventory, Supply Chain, and Transportation Systems

by

## Guillaume Roels

## Abstract

This thesis investigates the impact of lack of information and decentralization of decision-making on the performance of inventory, supply chain, and transportation systems. In the first part of the thesis, we study two extensions of a classic single-item, single-period inventory control problem: the "newsvendor problem." We first analyze the newsvendor problem when the demand distribution is only partially specified by some moments and shape parameters. We determine order quantities that are robust, in the sense that they minimize the newsvendor's maximum regret about not acting optimally, and we compute the maximum value of additional information. The minimax regret approach is scalable to solve large practical problems, such as those arising in network revenue management, since it combines an efficient solution procedure with very modest data requirements. We then analyze the newsvendor problem when the inventory decision-making is decentralized. In supply chains, inventory decisions often result from complex negotiations among supply partners and might therefore lead to a loss of efficiency (in terms of profit loss). We quantify the loss of efficiency of decentralized supply chains that use price-only contracts under the following configurations: series, assembly, competitive procurement, and competitive distribution.

In the second part of the thesis, we characterize the dynamic nature of traffic equilibria in a transportation network. Using the theory of kinematic waves, we derive an analytical model for traffic delays capturing the first-order traffic dynamics and the impact of shock waves. We then incorporate the travel-time model within a dynamic user equilibrium setting and illustrate how the model applies to solve a large network assignment problem.

# Acknowledgments

First and foremost, I would like to thank my advisor, Georgia Perakis, for her guidance throughout my PhD. It has been a great pleasure to work with her. Her diversity of research interests and her intellectual dynamism have enlightened my experience at MIT. I have really appreciated her commitment to my professional fulfillment, by encouraging me to submit papers for publication, to attend conferences, and to participate to competitions. In particular, Georgia has been extremely supportive during my job search, by living it as intensively as I did. Her mentorship will inspire me throughout my academic career. Her notorious dedication to the well-being of her students helped me to go through the ups and downs of my graduate student life, and I am thankful for her boundless support and friendship.

I would also like to thank Steve Graves and David Simchi-Levi for their comments and advices that have significantly improved the quality of this work and have opened directions for future research. I have been impressed by their ability to combine theory with practice and to deliver clear managerial insights, both in research and in teaching. This thesis is also based on the extensive work of Dimitris Bertsimas. I thank him for his challenging questions and his encouragements that comforted me in my choices of research topics and career.

I am grateful to Yves Pochet, Philippe Chevalier, and Laurence Wolsey for encouraging me to pursue a doctoral degree at MIT. They initiated me to operations research, made me discover the academic world, and have been great mentors since then. Yves detected early on my interest in research and developed my technical abilities to do research.

The parties, dinners, and expeditions with my ORC friends Mike, Margret, Katy, Jose, Victor, Felipe, Yann, Tim, and the visitors, Jeff, Samuel, and Jean-Philippe, are unforgettable. I also thank Theo, Ilan, Kostas, Ruben, Juliane, Mike M., Elodie, Carine, Carol, Shobhit, Pranava, Rags, and Nico for being such great colleagues, and Paulette, Laura, and Andrew for making the center run so smoothly.

I am also grateful to my parents, brothers and sisters, family in law, and friends

5

for supporting my dream to live abroad. I have enjoyed welcoming them in our small apartment and exploring with them the hidden gems of New England.

Last but not least, I would like to express my love to my wife, Charlotte. With her, the destination has been the journey. Our adventure abroad has been source of mutual care, tenderness, and love. She also gave me the most wonderful gifts: Elliot, whose smiles illuminate every day, and a second baby who will be no less charming.

# Contents

## Conclusions               161

## A  Proofs               163

# List of Figures

12

# List of Tables

# Structure of the Thesis

The thesis investigates the impact of lack of information and decentralization of decision-making on the performance of inventory, supply chain, and transportation systems. The first part of the thesis analyzes two extensions of the traditional "newsvendor problem," which is a single-item, single-period inventory control problem.

The newsvendor model, like most stochastic inventory models, assumes full knowledge about the demand probability distribution. However, in practice, it is often difficult to completely characterize the demand distribution, especially in fast-changing markets. In Chapter 1, we study the newsvendor problem with partial information about the demand distribution (e.g., mean, variance, symmetry, and unimodality). Specifically, we derive the order quantities that minimize the newsvendor's maximum regret about not acting optimally. Most of our solutions are tractable, which makes them attractive for practical application. We also quantify the "Price of Information," defined as the maximum opportunity cost from not knowing accurately the demand distribution. We then illustrate how the minimax regret criterion can be adopted to solve complex inventory problems, such as those arising in network revenue management, since our approach combines an efficient solution procedure with modest data requirements.

In Chapter 2, we quantify the loss of efficiency (i.e., loss of profit) associated with decentralizing inventory decision-making. In particular, we measure the efficiency of decentralized supply chains that use price-only contracts. With a price-only contract, a buyer and a seller agree only on a constant transaction price, without specifying the amount that will be transferred. It is well known that these contracts do not

provide incentives to the parties to coordinate their inventory/capacity decisions. Efficiency is measured using the "Price of Anarchy" ratio, defined as the largest ratio of profits between the integrated supply chain and the decentralized supply chain. We characterize the efficiency of various supply chain configurations: push or pull inventory positioning, two or more stages, serial or assembly systems, single or multiple competing suppliers, and single or multiple competing retailers.

In the second part of the thesis, we analyze a model of traffic assignment, called the "dynamic user equilibrium," in which travelers noncooperatively seek to minimize their travel time. To determine such an equilibrium, it is important to evaluate the impact of congestion on a driver's travel time. In Chapter 3, we derive an analytical travel-time function, based on the theory of kinematic waves, and integrate it within a dynamic user equilibrium setting. Our travel-time model captures the first-order traffic dynamics as well as the impact of shock waves. Numerical examples demonstrate the quality of the analytical function, in comparison with simulated travel times. We also prove that the travel-time function is continuous and strictly monotone if the flow varies smoothly. We then illustrate how the model can be used to solve a large network assignment problem through a numerical example.

All proofs appear in the appendix.

# Part I

# Beyond the Newsvendor Model

# Introduction to Part I

The newsvendor model is a single-item, single-period stochastic inventory model. A make-to-stock firm has to decide how much to order before a selling season, without knowing the demand, so as to maximize its expected profit.

The decision variable is the order quantity, denoted by $y$. Demand $D$ is random and is assumed to follow a distribution $F(x)$. One cannot sell more than what has been ordered; therefore, the amount sold will be $\min\{y, D\}$. We assume that the firm has a unit ordering cost $c$ and a unit selling price $p > c$. Units purchased but not sold become obsolete by the end of the season and have no salvage value. If the firm orders $y$ units, it pays an order cost $cy$ and receives an expected revenue equal to $pE[\min\{y, D\}]$. Therefore, the newsvendor problem can be formulated as follows:

$$\max_{y} E[\Pi(y, D)], \tag{1}$$

where $E[\Pi(y, D)] = pE[\min\{y, D\}] - cy$.

The name of the model is derived from the situation of a newsvendor who has to determine how many newspapers to buy at the beginning of the day, before attempting to sell them at a street corner. Despite the name that focuses on a specific application, the newsvendor model has been used in many different contexts, such as fashion retailing (Fisher and Raman 1996, Graves and Parsons 2005), capacity planning (Van Mieghem and Rudi 2002), and airline revenue management (Littlewood 1972).

The problem (1) is a concave optimization problem. The optimal order quantity $y^*$ is the smallest $y$ such that $F(y) \geq 1 - r$, where $r \doteq c/p$. With a nonnegative demand distribution, the optimal order quantity $y^*$ is also nonnegative. If the demand

19

distribution is continuous, the optimality condition is equivalent to

$$F(y^*) = 1 - r,$$

and the optimal order quantity is called the *critical fractile* solution, first appearing in Arrow, Harris, and Marschak (1951). Arrow (1958) attributes the model itself to Edgeworth (1888). The newsvendor model has been used as a building block for more complex inventory control problems including multiple periods, multiple items, multiple stages, fixed order costs, and pricing decisions among others; see Porteus (1990) for a review.

In this first part of the thesis, we investigate two extensions of the newsvendor model. First, the newsvendor model assumes full knowledge of the demand probability distribution; however, inventory decisions must often be made with little information about demand. Second, the newsvendor model assumes a single decision-maker; however, supply chains are composed of several decision-makers, possibly with conflicting objectives. In the thesis, we quantify how the lack of information about the demand distribution and the decentralization of decision-making affect the order quantity and the efficiency of the decision.

**Lack of information.** The newsvendor model (1) assumes full knowledge of the demand probability distribution. However, demand is often subject to several factors that are beyond the firm's control (competitors' prices, availability of alternative products, consumption behavior) and that are in essence hard to predict (Huyett 2005). Moreover, as economic conditions are changing fast (because of increased globalization, economic liberalization, etc.), companies introduce new products with very little historical data (Fisher and Raman 1996). Finally, many firms forecast a single point of demand using qualitative methods (Darlymple 1988), without estimating the full distribution.

In order to use the newsvendor model (1), one must select a demand distribution as an input to the model. But which distribution should be picked? Uniform, normal,

gamma, or exponential? All these distributions give rise to different order quantities. Ideally, the order quantity must be robust, i.e., perform well under most demand scenarios.

In Chapter 1, we examine the newsvendor model (1) with partial information about the demand distribution (e.g., range, mean, variance, symmetry, and unimodality). We derive the order quantities that minimize the newsvendor's maximum opportunity cost, called the *minimax regret*, from not knowing the demand distribution. We then show how the minimax regret approach can be used for solving practical problems. In particular, we focus on a network revenue management application, and illustrate how the minimax regret approach relies on modest data requirements and leads to efficient solution techniques.

**Decentralization of decisions.** The newsvendor model (1) assumes a single decision-maker. However, different partners in a supply chain have different exposures to risk; as a result, the level of inventory in a decentralized supply chain might be different from that in an integrated supply chain. In fact, it has been shown that price-only contracts, i.e., contracts that only specify the wholesale-price between two parties, do not coordinate the inventory decisions in a supply chain and lead to a profit loss; see Lariviere and Porteus (2001) and Cachon and Lariviere (2001).

To remedy this loss of efficiency, many more elaborate contracts have been proposed: buyback, revenue-sharing, nonlinear pricing schemes, etc. (see Cachon 2003 for a review). However, these contracts are in general more costly to negotiate and more complex to administrate. On the other hand, they have the advantage of coordinating the inventory decisions in a supply chain—at least under certain conditions.

Before implementing these more elaborate contracts, it is worth measuring the inefficiency associated with the simplest ones. In Chapter 2, we quantify the loss of efficiency associated with price-only contracts in decentralized supply chains. We measure the loss of efficiency using the "Price of Anarchy" ratio, defined as the largest ratio of profits between the integrated supply chain and the decentralized supply chain. We characterize the efficiency of the following supply chain configurations:

push or pull inventory positioning, two or more stages, serial or assembly systems, single or multiple competing suppliers, and single or multiple competing retailers.

**Notations.** We begin by introducing some notational conventions. Vectors are denoted in small bold letters. For a vector $\mathbf{x}$, $x_j$ denotes its $j$th component and $x_{-j}$ denotes the other components. All vectors are column vectors, and $\mathbf{x}'$ is the vector transpose. The function $\min\{\mathbf{x}, \mathbf{y}\}$ takes the componentwise minimum of vectors $\mathbf{x}$ and $\mathbf{y}$. Similarly, the function $\mathbf{x}^+$ takes the componentwise maximum of $\mathbf{x}$ and $\mathbf{0}$. Let $\mathbf{1}$ be a vector of ones. Matrices are denoted in capital bold letters. For a matrix $\mathbf{A}$, let $\mathbf{A}_j$ represent its $j$th column

Random variables are usually denoted with a capital letter while their realization is noted with a small letter; for example, $d$ is a realization of random variable $D$. A cumulative probability distribution $F(x)$ has a complementary distribution $\bar{F}(x) = 1 - F(x)$. If the distribution is continuous, it has a density $f(x)$. Finally, w.p.1 is short for *with probability 1*.

# Chapter 1

# The Price of Information in the Newsvendor Model

## 1.1 Introduction

The classic newsvendor model (1) assumes full knowledge of the demand probability distribution. However, in practice, it is often difficult to completely characterize the demand distribution, especially in fast-changing markets.

In this chapter, we study the newsvendor problem (1) with partial information about the demand distribution (e.g., mean, variance, symmetry, and unimodality). In particular, we derive the order quantities that minimize the newsvendor's maximum regret about not acting optimally. Most of our solutions are tractable, which makes them attractive for practical application. We also quantify the "Price of Information," defined as the maximum opportunity cost from not knowing accurately the demand distribution, and highlight that the shape of the demand distribution has sometimes more informational value than its variance. Finally, we illustrate the potential of the minimax regret approach for solving practical stochastic inventory control problems, such as those arising in network revenue management.

Decision-making models can be classified according to the degree of uncertainty they are addressing (Luce and Raiffa 1957). (a) If the firm orders after observing the demand realization, its decision is made *under certainty*. (b) If the firm orders before

the realization of the demand, but with knowledge of the demand distribution, its decision is made *under risk*. The classic newsvendor model (1) fits into this class of problems. (c) Finally, if the firm has no knowledge about the demand distribution, its decision is made *under uncertainty* or *under ambiguity*.

In this chapter, we adopt a median approach between risk and uncertainty, and assume only partial information about the demand distribution (e.g., range, mean, variance, symmetry, and unimodality). Partial demand information can come from expert judgments (i.e., depending on what the sales managers are confident in assuming about the demand) or a stationary forecast methodology (e.g., if the mean and variance of forecast errors are stationary over time, but the distribution of forecast errors is not). The problem of making an order decision with only partial information about the demand distribution was raised at the origins of inventory theory (Scarf 1958). In the next section, we review the different approaches that have been proposed subsequently to address this issue.

In a multi-period setting, past sales data are available and the optimal order quantity can be computed either with a Bayesian approach (e.g., Scarf 1960, Lariviere and Porteus 1999, Ding, Puterman, and Bisi 2002), or with a nonparametric approach (e.g., Godfrey and Powell 2001, Levi, Roundy, and Schmoys 2005). In contrast, we focus on the single-period nature of the game and assume that past sales data are unavailable.

When demand can be only vaguely described (e.g., "demand is *about d*"), fuzzy set theory can be applied to derived order quantities (see e.g., Petrović et al. 1996 and Guo 2003). In contrast, we assume precise, but partial, information about the demand distribution.

This chapter analyzes the newsvendor problem with partial information about the demand distribution. The main contributions of this chapter are the following:

1. We adopt an approach to solve a variety of problems that require a robust but not conservative solution.

2. We characterize the robust order quantities in the presence of partial infor-

mation about demand, such as the moments (mean, variance) and the shape (range, symmetry, and unimodality) of the demand distribution.

3. We compute the maximum value of additional information about the demand distribution, called the "Price of Information." In particular, we observe that the shape of the demand distribution can have more informational value than its variance.

4. We justify the use of entropy-maximizing distributions in the newsvendor model.

5. We illustrate how the minimax regret approach can be used to solve more complex problems, such as those arising in network revenue management.

The chapter is organized as follows. Section 1.2 reviews the common criteria for making decisions under uncertainty and motivates the minimax regret as a less conservative approach. In Section 1.3, we characterize the problem of minimizing the maximum regret and formulate it as a moment bound problem. In Section 1.4, we derive the order quantities that minimize the newsvendor's maximum regret. In Section 1.5, we quantify the Price of Information, and show that knowing the variance is in general less valuable than characterizing the shape of the demand distribution. In Section 1.6, we analyze a multi-item generalization of (1) with capacity constraints, and illustrate how the minimax regret approach can be used to solve complex network revenue management problems. Finally, we outline our conclusions in Section 1.7.

## 1.2 Decision-Making under Uncertainty

In this section, we review the common criteria for making decisions under uncertainty. Instead of assuming a particular probability distribution $F(x)$ in (1), we assume that the demand distribution is only partially specified (by its range, moments, shape, etc.). Let $\mathcal{D}$ be the set of probability distributions that are consistent with the prior information. Slightly abusing of notations, we write $D \in \mathcal{D}$ to say that the probability distribution of $D$ belongs to $\mathcal{D}$.

## 1.2.1 Maximin

The traditional paradigm to decision-making under uncertainty is the maximin approach and consists in maximizing the worst-case profit, where the worst case is taken over all possible distributions in $\mathcal{D}$. That is,

$$\max_{y} \min_{D \in \mathcal{D}} E[\Pi(y, D)]. \tag{1.1}$$

Scarf (1958) derived the optimal order quantity when only the mean $\mu$ and the variance $\sigma^2$ of the demand are known. Gallego and Moon (1993) rederived Scarf's result using a simpler proof and extended the model to include recourse orders, fixed order costs, random yields, and multiple products. The maximin approach has also been applied to multi-period inventory models, under continuous or periodic review (Gallego 1992, Moon and Gallego 1994, Gallego 1998), to the newsvendor model with customer balking (Moon and Choi 1995), to the newsvendor model with combined choice of lead-time and order quantity (Ouyang and Wu 1998), and to finite-horizon models with discrete demand distributions, incompletely specified by selected moments and percentiles (Gallego, Ryan, and Simchi-Levi 2001). Zhou and Natarajan (2005) also derived the minimax order quantity when the demand distribution is symmetric, with known mean and variance.

By definition, the maximin approach is conservative since it focuses on the worst-case profit. In some situations (e.g., when only the mean is known), the maximin criterion recommends not to order at all. In fact, the maximin approach guarantees a lower bound on the expected profit but says nothing about the quality of the solution under different demand scenarios.

To reduce the level of conservativeness of the approach, Ben-Tal and Nemirovski (1999) and Bertsimas and Sim (2004) suggested to reduce the size of $\mathcal{D}$ by imposing a "budget of uncertainty." By varying the budget, the decision-maker can adjust the level of conservativeness of the model. In particular, when the budget of uncertainty is zero, the set $\mathcal{D}$ is a singleton and (1.1) reduces to (1). Their approach performs well with a large number of random variables, as illustrated by Bertsimas and Thiele

(2006) in a multi-period environment, and we explore it in more detail in a multi-item extension of (1) in Section 1.6.

There are other criteria for decision-making under uncertainty that are similar to the maximin (Luce and Raiffa 1957): the maximax, often considered as too optimistic; and the Hurwicz criterion, which combines the maximin and the maximax criteria.

## 1.2.2 Minimax Regret

The concept of regret was introduced by Savage (1951), as an improvement over the too-conservative maximin criterion. For a particular random event $D$, the expected profit associated with a decision $y$ equals $E[\Pi(y, D)]$. If the probability distribution of $D$ was known, the maximum expected profit would be equal to $\max_z E[\Pi(z, D)]$. The regret measures the additional expected profit that could have been obtained if the probability distribution had been known at the time of the decision, that is

$$\rho(y, D) = \max_z E[\Pi(z, D)] - E[\Pi(y, D)]$$

While decision $y$ is distribution-free, decision $z$ depends on the probability distribution. Similarly, the relative regret, analog to the "competitive ratio" (e.g., see Karp 1992), measures the additional profit increase that could have been obtained with full information, and is defined as $\max_z E[\Pi(z, D)]/E[\Pi(y, D)]$.

Savage (1951) proposed to minimize the maximum regret as a decision criterion under uncertainty, where the maximum is taken over all possible probability distributions in $\mathcal{D}$, that is

$$\min_y \max_{D \in \mathcal{D}} \rho(y, D). \tag{1.2}$$

The minimax regret was introduced as an improvement over the maximin criterion. While the maximin criterion is pessimistic, the minimax regret is more informed, as it benchmarks any decision against the optimum under full information. Among the distributions in $\mathcal{D}$, some lead to more extreme order decisions than others. In general,

27

the extreme order quantities (e.g., such as the solution to the maximin objective) are associated with distributions at the boundary of $\mathcal{D}$. The purpose of the minimax regret approach is to get away from the boundaries of $\mathcal{D}$ and to choose a decision that performs well under most probability distributions.

The regret is a pure intellectual construct. In practice, there does not exist a "true" probability distribution, especially in a single-period setting, and the actual opportunity cost will never be measured. Instead, the set of probability distributions $\mathcal{D}$ should be considered as expressions of the decision-maker's ignorance and is essentially defined subjectively (Keynes 1921, Savage 1954).

Savage (1951) introduced the concept of regret as a normative criterion, for making decisions under uncertainty. By contrast, Bell (1982) and Loomes and Sugden (1982) used the regret in a descriptive setting, in order to justify some behaviors that the classic expected utility theory considers as paradoxical (e.g., the coexistence of gambling and insurance). In the context of the newsvendor model, Schweitzer and Cachon (2000) and Brown and Tang (2006) experimentally observed a systematic departure from the traditional newsvendor solution. In fact, managers tend to make order decisions that minimize the *ex-post* difference between inventory and demand realization, independently of the costs involved.

The concept of regret as a normative criterion for decision-making (1.2) has been applied to different contexts, with different specifications of $\mathcal{D}$. Chamberlain (2000) considered a parametric family of distributions. Bergeman and Schlag (2005) and Lim and Shanthikumar (2006) defined $\mathcal{D}$ as a neighborhood around a distribution of reference. When $\mathcal{D}$ is defined as the set of distributions with a given support, the regret approach has been applied to solving combinatorial problems with uncertain cost parameters (Kouvelis and Yu 1997), with a recent attention on the knapsack problem (Averbakh 2001, Conde 2004) and on linear optimization problems (Inuiguchi and Sakawa 1995, Averbakh and Lebedev 2005).

In inventory management, the regret has been used by Morris (1959), Kasugai and Kasegai (1961), and Vairaktarakis (2000) in the newsvendor model, and Yu (1997) in the EOQ model, when only the support of the demand is known. Yue, Chen, and

Wang (2006) used the regret in the newsvendor model when the mean and the variance of the demand are known, with no restriction on the nonnegativity of demand.

In this chapter, we define $\mathcal{D}$ as the set of distributions with certain moments and shape, similarly to Scarf (1958) and Gallego and Moon (1993) with the maximin criterion, and Yue, Chen, and Wang (2006) with the minimax regret criterion.

In addition to providing a robust decision, solving problem (1.2) quantifies the maximum opportunity cost from not having full information about the demand distribution. If the newsvendor has the option to conduct a marketing survey to learn more about the demand distribution, she will evaluate this option against its potential gain, which suggests the next definition.

**Definition 1.** *The Price of Information (PoI) corresponds to the maximum profit loss from knowing only partial information about the demand distribution. Mathematically, the PoI is the optimal value of (1.2).*

## 1.2.3 Entropy Maximization

In contrast to the maximin and minimax regret, entropy maximization is not a decision criterion, but is a criterion for selecting a probability distribution among a set $\mathcal{D}$, as an input to a stochastic decision model. The principle of insufficient reason, proposed by Laplace, states that, with no information available, all possible outcomes should be considered as equally likely (Luce and Raiffa 1957). Jaynes (1957, 2003) extended the principle of insufficient reason by proposing to consider the distribution that maximizes the entropy over the set of distributions $\mathcal{D}$.

For a discrete random variable, with probabilities $P_1, ..., P_n$, the entropy, as defined by Shannon in the field of information theory, is equal to

$$H(P_1, ..., P_n) = -K \sum_i P_i \ln P_i,$$

where $K$ is a positive constant. The entropy of a probability distribution represents the amount of uncertainty associated with the distribution. Jaynes (1957) claimed that the distribution that maximizes the entropy is a good prior distribution, as

29

it is the "maximally noncommittal with regard to missing information." In fact, the entropy is similar to the barrier function in nonlinear programming (Bertsekas 1999); the distribution that maximizes the entropy is therefore equivalent to the analytical center of $\mathcal{D}$. The entropy has also been defined for continuous distributions, in reference to a prior measure (usually taken as uniform); see Jaynes (2003) for details.

Table 1.1 lists some of the distributions that maximize the entropy for a given prior information. The distribution that maximizes the entropy among all distributions with given mean and variance (without restrictions on nonnegativity) is the normal distribution. While the Central Limit Theorem explains the ubiquity of the normal distribution, the maximum-entropy principle justifies the assumption of normal distribution when the first two moments of a random variable are specified.

Table 1.1: Entropy-maximizing distributions

| Prior information | MaxEnt Distribution |
|---|---|
| Range | Uniform |
| Mean and nonnegativity | Exponential |
| Mean and variance | Normal |

Since the selection of the demand distribution is distinct from the costs of the decision (Hayes 1969), there is no guarantee that maximizing entropy is a reasonable criterion for choosing a distribution as an input to (1). In contrast, the minimax regret approach addresses both problems (selection of the demand distribution and choice of order quantity) simultaneously and has a performance guarantee.

## 1.3   Methodology

In this section, we formulate the minimax regret problem (1.2) as a moment bound problem. Moment bound problems aim at maximizing a function over all possible random distributions that satisfy some moment constraints. They have been applied to deriving Chebyshev-type inequalities, pricing options, and managing inventories (Scarf's maximin problem can also be formulated as a moment bound problem).

Suppose that the newsvendor considers a convex class of distributions whose $n$ first moments match $\mu$, denoted by $\mathcal{D}$. For convenience, we take $\mu_0 = 1$. If $n$ moments are known, any distribution $F \in \mathcal{D}$ satisfies the following constraints:

$$\int_0^\infty x^i dF(x) = \mu_i, \ \forall i = 0, ..., n.$$

We also assume that certain Slater conditions hold on the moment constraints (i.e., the moment vector is interior to the set of feasible moments).

Problem (1.2) can be reformulated as follows, by inverting the order of maximization:

$$
\begin{aligned}
& \min_y \max_{D \in \mathcal{D}} \max_z E[\Pi(z, D) - \Pi(y, D)] \\
= \ & \min_y \max_z \left\{ \max_{D \in \mathcal{D}} E[\Pi(z, D) - \Pi(y, D)] \right\}, \\
= \ & \min_y \max_z p \left\{ \max_{F \in \mathcal{D}} \int_0^\infty (\min\{x, z\} - \min\{x, y\}) dF(x) \right\} + c(y - z). \quad (1.3)
\end{aligned}
$$

Let $\mathcal{D}$ be the convex class of distributions that satisfy the shape constraints, and $\Omega$ be the support of the distribution. The inner problem can then be formulated as the following moment bound problem:

$$
\begin{aligned}
\max_{F \in \mathcal{D}} \quad & \int_\Omega (\min\{x, z\} - \min\{x, y\}) dF(x), \\
\text{s.t.} \quad & \int_\Omega x^i dF(x) = \mu_i, \ \forall i = 0, ..., n.
\end{aligned}
\quad (1.4)
$$

If $\mathcal{D}$ is simply the set of nonnegative distributions with support $\Omega$, the above problem is a semi-infinite linear optimization problem, in which the variables are the amount of mass at each point $x \in \Omega$. By the theory of linear optimization, at most $n + 1$ variables are positive in an optimal solution. Hence, the distribution achieving the maximum regret is a discrete distribution, with mass at $n + 1$ points at most (Smith 1995, Bertsimas and Popescu 2002, 2005).

For a general convex set of probability measures $\mathcal{D}$, problem (1.4) is equivalent to its relaxation over the cone of measures, i.e., when $F$ belongs to the cone generated by the set $\mathcal{D}$. By strong duality (under Slater's conditions), problem (1.4) is equivalent

to the following dual problem (Popescu 2005):

$$\min_{\alpha_0,\ldots,\alpha_n} \quad \sum_{i=0}^{n} \alpha_i \mu_i,$$
$$\text{s.t.} \quad \sum_{i=0}^{n} \alpha_i x^i - (\min\{x,z\} - \min\{x,y\}) \in \mathcal{C}^*, \tag{1.5}$$

where $\mathcal{C}^*$ is the polar of the cone of measures defined by $\mathcal{D}$. When $\mathcal{D}$ is the set of nonnegative distributions with support $\Omega$, the dual problem (1.5) simplifies to the following semi-infinite linear optimization problem:

$$\min_{\alpha_1,\ldots,\alpha_0} \quad \sum_{i=0}^{n} \alpha_i \mu_i,$$
$$\text{s.t.} \quad \sum_{i=0}^{n} \alpha_i x^i \geq \min\{x,z\} - \min\{x,y\}, \quad \forall x \in \Omega. \tag{1.6}$$

Bertsimas and Popescu (2002, 2005) showed that the above semi-infinite linear optimization problem can be formulated as a semi-definite optimization problem and be therefore efficiently solved. Similar results hold for the general conic dual problem, under certain conditions (Popescu 2005).

The next proposition characterizes the optimal solution to problem (1.3).

**Proposition 1.** *(a) The function* $\Phi(z;y) = \max_{D \in \mathcal{D}} E[\Pi(z,D) - \Pi(y,D)]$ *is quasi-concave on the interval* $z \in [0,y]$ *and on the semi-interval* $z \in [y,\infty)$, *but not necessarily on* $[0,\infty)$.

*(b) The function* $\Upsilon(y) = \max_z \{\max_{D \in \mathcal{D}} E[\Pi(z,D) - \Pi(y,D)]\}$ *is a convex function of* $y$.

Part (a) states that the optimization problem over $z$ is not a concave problem, unless we distinguish the two cases $y \leq z$ and $y \geq z$. According to Part (b), the optimal robust quantity $y$ equates the regret from ordering too little to the regret from ordering too much (minimum of two convex functions).

## 1.4 Minimax Regret Order Quantities

In this section, we derive the order quantities that minimize the newsvendor's maximum regret, in the presence of limited information about demand.

**Range.** When the demand distribution is known to lie in some bounded range, the maximum regret order quantity is given by a convex combination of the boundaries of the interval, as stated in the following theorem.

**Theorem 1.** *If the demand distribution is nonnegative, with support $[l, u]$, the minimax regret order quantity is the following:*

$$y = rl + (1 - r)u,$$

*where* $r \doteq c/p$, *and the Price of Information amounts to*

$$PoI = c(1 - r)(u - l).$$

The minimax regret order quantity is the same as the optimal solution of the newsvendor model (1) with a uniform demand distribution, consistently with the principle of insufficient reason (or the maximum entropy), which suggests a uniform prior when only the range is known (Table 1.1).

**Range and Mean.** Suppose that the newsvendor knows the mean of the demand. This case corresponds for instance to the situation in which a firm has a single-point forecast of demand. We first derive a general result, if the demand distribution is known to have its support on $[l, u)$, and specialize it after to the case where $l = 0$ and $u = \infty$.

**Theorem 2.** *If the demand distribution is nonnegative, with support $[l, u)$ and mean* $\mu$, *the minimax regret order quantity is the following:*

$$y = \begin{cases} lr + \mu(1 - r), & \text{if } \frac{1}{2} \leq r, \\ l + \frac{\mu - l}{4r}, & \text{if } \frac{1}{2}\frac{\mu - l}{u - l} \leq r \leq \frac{1}{2}, \\ u - r\frac{(u - l)^2}{\mu - l}, & \text{if } r \leq \frac{1}{2}\frac{\mu - l}{u - l}, \end{cases}$$

*and the Price of Information amounts to*

$$PoI = \begin{cases} c(1-r)(\mu - l), & \text{if } \frac{1}{2} \le r, \\ c\frac{\mu-l}{4r}, & \text{if } \frac{1}{2}\frac{\mu-l}{u-l} \le r \le \frac{1}{2}, \\ c(u-l)(1 - r\frac{u-l}{\mu-l}), & \text{if } r \le \frac{1}{2}\frac{\mu-l}{u-l}. \end{cases}$$

The support of the demand distribution plays a similar role to the "budget of uncertainty" proposed by Bertsimas and Sim (2004). By restricting the size of the support of the demand distribution, the decision-maker adjusts the amount of variability she wants to cover with her decision. Smaller intervals correspond to higher degrees of confidence in the average value of the demand.

**Mean.** The next corollary derives the optimal order quantity when the newsvendor knows nothing about the support of the demand distribution, i.e., when $l = 0$ and $u = \infty$.

**Corollary 1.** *If the demand distribution is known to be nonnegative with mean $\mu$, the minimax regret order quantity is equal to*

$$y = \begin{cases} \mu(1-r), & \text{if } \frac{1}{2} \le r, \\ \frac{\mu}{4r}, & \text{if } \frac{1}{2} \ge r, \end{cases} \tag{1.7}$$

*and the Price of Information amounts to*

$$PoI = \begin{cases} c(1-r)\mu, & \text{if } \frac{1}{2} \le r, \\ c\frac{\mu}{4r}, & \text{if } \frac{1}{2} \ge r. \end{cases}$$

When only the mean is known, the maximin approach recommends not to order since the worst-case demand scenario is close to a unit impulse at zero. In contrast, the minimax regret approach balances the losses incurred with low demand scenarios with the opportunity costs associated with high demand scenarios.

34

The exponential distribution maximizes the entropy over all nonnegative distributions with mean $\mu$ (Table 1.1). The optimal order quantity of (1) with an exponential distribution is equal to $-\ln(r)\mu$. Figure 1-1 compares the minimax regret order quantity with $-\ln(r)\mu$ for increasing values of $1/r$, as well as the expected profits when the demand distribution is exponential. Although the minimax regret order quantities are lower than the order quantities derived with the exponential distribution, especially for large profit margins, the expected profits are similar when demand is exponentially distributed.



Figure 1-1: Order quantities and expected profits when only the mean of the demand distribution is known ($\mu = 100$)

Table 1.2 compares the average performance of the minimax regret approach to the traditional newsvendor approach, with a gamma distribution. The gamma distribution has two parameters: the scale $\theta$ and the shape $k$. Its mean equals $k\theta$ and its coefficient of variation equals $1/\sqrt{k}$. We consider three demand scenarios, each with a mean 100, and with $k$ equal to .1, 1, and 10. Table 1.2 shows the expected profits when $p = 1.5$, $c = 1$ (left part) and $p = 3$, $c = 1$ (right part). For each demand scenario (i.e., value of $k$), we compare the order quantities $y$ and the expected profit obtained with the minimax regret approach and the solution of (1). The maximum

expected revenues correspond to the diagonal (when the correct value of $k$ is used in the newsvendor model).

When the profit margins are small ($p = 1.5$, $c = 1$), the potential losses from choosing an erroneous parameter mostly result from holding too much inventory. The robust policy avoids most of the losses when $k$ is small, but only captures 50% of the profit when $k$ is large. On the other hand, when the profit margins are higher ($p = 3$, $c = 1$), the losses from choosing a wrong parameter are essentially opportunity costs from not meeting the demand. The robust policy avoids most of the losses when $k$ is small and captures more than 85% of the profit when $k$ is large.

Table 1.2: Order quantities and expected profits with $\mu = 100$

|  | $p = 1.5, c = 1$ | | | | $p = 3, c = 1$ | | | |
|---|---|---|---|---|---|---|---|---|
| Approach | $y$ | $k = .1$ | $k = 1$ | $k = 10$ | $y$ | $k = .1$ | $k = 1$ | $k = 10$ |
| Minimax regret | 33.33 | -17.28 | 9.18 | 16.65 | 75 | -15.36 | 83.29 | 141.02 |
| $\Gamma$ with $k = 0.1$ | 0.01 | 0.00 | 0.01 | 0.01 | 10.63 | 1.92 | 19.63 | 21.27 |
| $\Gamma$ with $k = 1$ | 40.54 | -21.82 | 9.45 | 20.20 | 109.86 | -31.51 | 90.14 | 164.40 |
| $\Gamma$ with $k = 10$ | 83.94 | -51.57 | 1.26 | 33.76 | 110.66 | -31.92 | 90.14 | 164.41 |

**Mean=Median.** When the mean equals the median of the distribution, there is 50% chance that the demand realization falls below the mean $\mu$ and 50% chance that it is above $\mu$. This happens when the demand distribution (or the distribution of forecast errors) is not skewed. In fact, most of the common distributions (e.g., normal, uniform, triangular) have the median corresponding to the mean. With this additional restriction, the robust order quantity changes substantially, as shown in the next theorem.

**Theorem 3.** *If the distribution is known to have its mean $\mu$ equal to the median, the minimax regret order quantity is equal to*

$$y = \begin{cases} 2\mu(1 - r), & \text{if } \frac{1}{4} \leq r, \\ \mu\frac{1+8r}{8r}, & \text{if } \frac{1}{4} \geq r, \end{cases}$$

36

*and the Price of Information amounts to*

$$
PoI = \begin{cases} c\mu(2 - \frac{1}{r})(1-r), & \textit{if } \frac{1}{2} \leq r, \\ c\mu(\frac{1}{r} - 2)r, & \textit{if } \frac{1}{4} \leq r \leq \frac{1}{2}, \\ c\frac{\mu}{8r}, & \textit{if } \frac{1}{4} \geq r. \end{cases}
$$

The additional requirement that the mean equals the median makes the order quantity larger than (1.7). In particular, it is now optimal to order more than the mean demand when $r$ falls below $1/2$, instead of $1/4$. Moreover, for low-margin products ($r \geq 1/2$), the newsvendor orders twice as much as (1.7).

**Mean and Symmetry.** Let us assume that the demand distribution is known to be symmetric with mean $\mu$. A symmetric demand distribution gives the same probability to $\mu - x$ and $\mu + x$, for all $0 \leq x \leq \mu$. Despite this case is stronger than having the mean equal to the median, the robust order quantity is the same as in Theorem 3, when $r \geq 1/4$.

**Theorem 4.** *When the demand distribution is known to be nonnegative, symmetric, with mean $\mu$, the minimax regret order quantity is equal to*

$$
y = 2\mu(1 - r),
$$

*and the Price of Information amounts to*

$$
PoI = \begin{cases} c\mu(2 - \frac{1}{r})(1-r), & \textit{if } 1/2 \leq r, \\ c\mu(\frac{1}{r} - 2)r, & \textit{if } 1/2 \geq r. \end{cases}
$$

From the nonnegativity of demand, a symmetric demand distribution with mean $\mu$ is bounded from above by $2\mu$. Interestingly, the optimal order quantity with a symmetric demand is the same as the optimal solution of the newsvendor model (1) with a uniform demand distribution. Intuitively, the uniform distribution is "in

the middle" of all symmetric distributions. At one extreme, we have the two-point distribution, giving positive probability to both 0 and $2\mu$; the other extreme is the deterministic case, with all probability mass at $\mu$.

The order quantity derived in Theorem 4 is the same as the one derived in Theorem 1, when the demand distribution is known to have its support on $[0, 2\mu]$. However, the robust value of additional information is lower in the case of symmetry.

**Mode.** Very often, the mode of the demand, i.e., the demand outcome that has the largest probability of occurring, is taken as a reference, instead of the mean demand. By focusing on the mode, the newsvendor exhibits some aversion against the ex-post difference between inventory and demand, similarly to the behavior observed by Schweitzer and Cachon (2000).

We assume that the demand distribution is unimodal, with its mode equal to $m$. A unimodal cumulative distribution function with mode $m$ is convex to the left of $m$ and concave to the right of $m$. If the distribution is continuous, the density function is increasing to the left of the mean and decreasing to the right.

When only the mode is known, the maximum regret is infinite because the regret when $z \geq y$, equal to $(p - c)(z - y)$, goes to infinity when $z$ increases. In fact, there is no upper bound on the critical fractile $z$ when only the mode is known. To avoid this problem, we assume in Theorem 5 that the demand distribution is nonnegative and bounded above by $u$.

**Theorem 5.** *If the demand distribution is known to be nonnegative, bounded from above by $u$, and unimodal with mode $m < u$, the minimax regret order quantity is*

$$
y = \begin{cases} \sqrt{m(1 - r)(2rm + u(1 - r))}, & \text{if } u \leq m(1 + (\frac{r}{1-r})^2), \\ u - \sqrt{r(u - m)(2u - ru + 2rm - 2m)}, & \text{if } u \geq m(1 + (\frac{r}{1-r})^2), \end{cases}
$$

38

*and the Price of Information amounts to*

$$PoI = \begin{cases} c\frac{1}{2}(\sqrt{(1-r)(2m+u(\frac{1}{r}-1))} - \sqrt{m}(\frac{1}{\sqrt{r}} - \sqrt{r}))^2, \\ \qquad\qquad\qquad\qquad\qquad if\ u \le m(1 + (\frac{r}{1-r})^2), \\ c(m(r/2-1) + u - \sqrt{r(u-m)(2u-ru+2rm-2m)}), \\ \qquad\qquad\qquad\qquad\qquad if\ u \ge m(1 + (\frac{r}{1-r})^2). \end{cases}$$

Figure 1-2 illustrates how the robust order quantity changes with the mode, for a fixed upper bound $u = 160$. When the mode is large, e.g., equal to 100, the minimax regret order quantity is close to the optimal solution of the newsvendor model with a normal demand distribution with a mean/mode equal to 100 and a standard deviation equal to 30 (so that the normal distribution is virtually bounded from above by $u$) or with a uniform demand distribution (for which the mode is anything in the interval $[0, 160]$). However, when the mode equals zero, as is common for slow-moving items which have a high probability of not being ordered but a fat right tail, the minimax regret order quantity is smaller.



Figure 1-2: Order quantities when the mode of the demand distribution is known

More general results can be derived when the mode is known to lie in a certain

39

interval $[m_1, m_2]$ (see Popescu 2005). However, solving the model analytically is cumbersome, as there are 12 cases to consider, depending on the relative order of $m_1, m_2, y$, and $z$.

**Mean, Unimodality, and Symmetry.** When the demand distribution is known to be symmetric, one of the worst-case demand scenarios is a two-point demand distribution, with nonzero probabilities at 0 and $2\mu$. Arguably, this worst-case scenario might not be realistic. With the additional requirement that the demand distribution is unimodal, this worst case is ruled out. As a result, the Robust Value of Information is reduced, as shown in the next theorem.

**Theorem 6.** *If the distribution is known to be nonnegative, symmetric, unimodal, with mean $\mu$, the minimax regret order quantity is equal to*

$$y = \begin{cases} 2\mu\sqrt{r(1-r)}, & \text{if } \frac{1}{2} \leq r, \\ 2\mu(1 - \sqrt{r(1-r)}), & \text{if } \frac{1}{2} \geq r, \end{cases}$$

*and the Price of Information amounts to*

$$PoI = \begin{cases} c\mu(1 - 2\sqrt{r(1-r)})\frac{1-r}{r}, & \text{if } \frac{1}{2} \leq r, \\ c\mu(1 - 2\sqrt{r(1-r)}), & \text{if } \frac{1}{2} \geq r. \end{cases}$$

Although the requirement of symmetry and unimodality might seem restrictive, many distributions fall into this category: uniform, normal, and Poisson if $\mu$ is large. Moreover, no assumption is made about the variance.

The right part of Figure 1-3 compares the regret minimax order quantities for various classes of demand distribution, assuming an average demand of 100. For comparison purposes, we have plotted on the left part of the figure the optimal solution to the newsvendor model (1), for a normal demand distribution with mean 100 and standard deviation of 1, 30, and 60.

As shown in the figure, the robust order quantity is increasing with $1/r$, under all scenarios. When only the mean of the demand is known, the order quantity is

Figure 1-3: Order quantities when the mean of the demand distribution is known in addition to some shape characteristics

lower than the average demand, for $1/r \leq 4$ (curve with triangles). When the demand distribution has more structure (mean/median, symmetry, or mode and upper bound), the optimal order quantity becomes closer to the order quantity obtained with the newsvendor model for a normal demand distribution. Moreover, the optimal order quantity is equal the mean demand when $r = 1/2$ (provided that $u = 2m$). Finally, when the demand distribution is both unimodal and symmetric with mean 100, the robust order quantity is almost equal to the mean, for sufficiently small $r$. Intuitively, the conditions of symmetry and unimodality tend to accumulate more probability mass about the mean, leading the order quantity to be also closer to the mean demand.

Comparing the left and right parts of Figure 1-3 illustrates the similarities of effect between restricting the shape of the demand distribution to a "regular" class and improving forecast accuracy (measured by a decrease of the standard deviation). As the standard deviation decreases, the optimal order quantities are also re-centered about the mean. Therefore, joint efforts could be devoted to both reducing forecast

error variability and better characterizing the shape of the demand distribution since both lead to less sensitive order quantities.

**Mean and Variance.** The variance of the demand distribution is usually taken as the variance of forecast errors. In contrast to the maximin approach, the minimax regret approach does not lead to a closed-form solution of the optimal order quantity.

**Theorem 7.** *If the demand distribution is known to be nonnegative, with mean $\mu$ and variance $\sigma^2$, the minimax regret order quantity is the solution to the following equality:*

$$
\max \left\{ \max_{\max\{\mu,y\} \leq x \leq \frac{\sigma^2+\mu^2}{\mu}} (p\frac{\mu}{x} - c)(x - y) \, , \right.
$$
$$
\left. \max_{y \leq x \leq \mu \cup \max\{y,\frac{\sigma^2+\mu^2}{\mu}\} \leq x \leq y + \sqrt{\sigma^2+(y-\mu)^2}} (p\frac{\sigma^2}{\sigma^2 + (x - \mu)^2} - c)(x - y) \right\} \quad (1.8)
$$
$$
= \left\{ \max_{\min\{\mu,y\} \geq x \geq \max\{0, y - \sqrt{\sigma^2+(\mu-y)^2}\}} (x - y)(p\frac{(x - \mu)^2}{\sigma^2 + (x - \mu)^2} - c) \right\}.
$$

When the variance grows to infinity, the minimax regret order quantity (1.8) tends to the robust order quantity when only the mean demand is known, derived in (1.7). At the other extreme, when the variance tends to zero, it becomes optimal to order exactly the mean demand.

Condition (1.8) involves three optimization problems. The first problem depends on the mean only, while the two other problems depend on both the mean and the variance. In fact, there is some similarity between the first problem and Markov's inequality and the next two problems and Chebyshev's inequality since the optimal solution of (1) is characterized by a critical fractile.

The robust order quantity can efficiently be found by a line search since, from Proposition 1, the inner maximization problems are quasi-concave. Alternatively, the robust order quantity can be approximated by considering a feasible solution, instead of an optimal solution, in the inner maximization problems of (1.8).

42

**Approximation 1.** *When $\sigma/\mu \leq \sqrt{1-r}$, the minimax regret order quantity, obtained by solving (1.8), can be approximated as follows:*

$$y \approx \max\{0, \mu + \sigma \frac{2}{5} \frac{1-2r}{\sqrt{r(1-r)}}\}. \tag{1.9}$$

The approximate order quantity (1.9) is the sum of the mean demand and a (possibly negative) safety stock. As in the traditional newsvendor model with a normal distribution, the safety stock is proportional to the standard deviation of the demand by a factor that is concave decreasing with $r$. The safety factor is negative for low-margin products ($r \geq 1/2$) and positive otherwise.

The approximation is constructed by considering a feasible solution, instead of the optimal solution, in the inner optimization problems in (1.8). In particular, we set $x = \mu + 1/2\sigma\sqrt{(p-c)/c}$ in the left-hand side, and $x = \mu - 1/2\sigma\sqrt{c/(p-c)}$ in the right-hand side. Moreover, we ignore the first maximization problem in the left-hand side. Equating both sides of the equality and solving for $y$, gives rise to approximation (1.9). More details are provided in the appendix.

The approximation is derived by equating the two optimization problems in (1.8) that depend on both the mean and the variance. In fact, the presence of the first optimization problem (related to Markov's inequality) is essentially due to the non-negativity requirement. If the distribution had its support on the real line, only the last two terms would appear in (1.8). Yue, Chen, and Wang (2006) analyzed the newsvendor problem with the minimax regret objective, when the mean and the variance are known, but with no restriction on the nonnegativity of demand. As a result, their safety stock factor (which is also not in closed-from) is almost equal to the approximated safety stock factor in (1.9), as it is illustrated in Table 1.3.

Table 1.3: Comparison between the safety factor in (1.9) and the safety factor obtained by ignoring the nonnegativity of the demand

| $1/r$ | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|
| $\frac{2}{5} \frac{1-2r}{\sqrt{r(1-r)}}$ | 0.0000 | 0.2828 | 0.4619 | 0.6000 | 0.7155 | 0.8165 |
| Yue et al. (2006) | 0.0000 | 0.2758 | 0.4503 | 0.5850 | 0.6977 | 0.7961 |

Table 1.4 compares the relative difference between approximation (1.9) and the minimax regret order quantity, obtained from (1.8), for various levels of profit margins and coefficients of variation. As expected, the quality of the approximation deteriorates when the coefficient of variation becomes large (as the probability of negative demand becomes non-negligible) and as $r$ increases.

Table 1.4: Relative difference between the robust order quantity and its approximation

| $1/r$ | $\sigma/\mu$ | | | | |
|-------|------|------|------|------|------|
| | .2 | .4 | .6 | .8 | 1 |
| 1.5 | 0.12% | 0.92% | 4.95% | 8.84% | 12.56% |
| 2 | 0.00% | 0.00% | 1.00% | 7.12% | 14.09% |
| 2.5 | 0.06% | 0.11% | 0.17% | 2.30% | 9.51% |
| 3 | 0.11% | 0.21% | 0.30% | 0.46% | 4.99% |

Figure 1-4 compares the order quantities obtained with the minimax regret objective (obtained by solving (1.8)), the maximin objective, and the newsvendor model (1) with a normal distribution. When the coefficient of variation is small, e.g., when



Figure 1-4: Order quantities when the mean and the variance of the distribution are known ($\mu = 100$)

$\mu = 100$ and $\sigma = 30$ as in the left part of Figure 1-4, there is almost no difference

among the three approaches. Similarly, Naddor (1978) observed that the optimal solution of (1) is insensitive to the choice of the demand distribution when the first two moments are known. As a result, selecting a normal demand distribution when the first two moments are known can be justified on the grounds of

- the Central Limit Theorem,

- maximizing the entropy,

- (approximately) maximizing the worst-case profit, and

- (approximately) minimizing the maximum regret.

However, when the coefficient of variation is large, assuming a normal distribution is no longer valid, as the probability of negative demand becomes non negligible. The right part of Figure 1-4 compares the order quantities when $\mu = \sigma = 100$. Since the coefficient of variation is large, the maximin approach recommends zero order quantities and the newsvendor solution recommends negative order quantities when the profit margin is small. In contrast, the minimax regret approach always recommends positive order quantities. Moreover, the order quantity obtained with the minimax regret objective is less sensitive to the profit margin, which is an attractive feature when costs are uncertain as well.

In addition, when the profit margin is large, and the coefficient of variation goes to infinity, the newsvendor order quantity obtained with a normal demand distribution grows to infinity, while the order quantities obtained with gamma, lognormal, and negative binomial distributions (which are nonnegative) tend to zero (Gallego, Katircioglu, and Ramachandran 2006). In fact, as the coefficient of variation increases, more probability mass is accumulated at zero for nonnegative distributions, while it is accumulated in the tails for a normal distribution. Gallego et al. also show that the order quantity associated with nonnegative demand distributions is bounded from above by $\mu/(1 - r)$, when the coefficient of variation becomes large, and that the bound is tight. Figure 1-5 compares the order quantity as the coefficient of variation increases, for $p = 20$, $c = 1$, using the newsvendor model with a normal

demand distribution, the maximin approach, and the minimax regret approach. The maximin approach recommends not to order at all if the coefficient of variation exceeds some threshold. In contrast, the minimax regret order quantity tend to $\mu/(4r)$ (see (1.7)) as the coefficient of variation becomes large, since having infinite variance is equivalent, from an informational perspective, to knowing only the mean.



Figure 1-5: Evolution of the order quantity when the coefficient of variation increases $(\mu = 100, p = 20, c = 1)$

Table 1.5 compares the average performance of the minimax regret approach with those of the maximin approach and the traditional newsvendor approach. We consider three demand scenarios, each with a mean equal to 100 and a standard deviation equal to 60: (truncated) normal (N), uniform (U), and gamma ($\Gamma$). When $p = 1.2$, $c = 1$, the expected profit with the maximin policy is zero. In contrast, the minimax regret quantity leads to a positive profit in all demand scenarios, and is in fact very close to the maximum achievable when the demand distribution is normal or gamma. On the other hand, when $p = 3$, $c = 1$, both robust approaches perform equally well, within a few percents of the maximum profit.

Table 1.5: Order quantities and expected profits with $\mu = 100$, $\sigma = 60$

| Approach | $p = 1.2, c = 1$ | | | | $p = 3, c = 1$ | | | |
|---|---|---|---|---|---|---|---|---|
| | $y$ | N | U | $\Gamma$ | $y$ | N | U | $\Gamma$ |
| Minimax regret | 56.97 | 1.40 | 1.10 | 5.44 | 116.62 | 133.76 | 129.58 | 132.11 |
| Maximin | 0 | 0 | 0 | 0 | 121.21 | 134.35 | 130.61 | 131.79 |
| N | 41.95 | 2.01 | 2.65 | 6.06 | 125.84 | 134.55 | 131.32 | 131.14 |
| U | 31.82 | 1.77 | 2.96 | 5.41 | 134.08 | 133.95 | 131.82 | 129.25 |
| $\Gamma$ | 44.84 | 1.99 | 2.46 | 6.09 | 114.54 | 133.36 | 129.02 | 132.14 |

## 1.5 Value of Additional Information

Using the mean and the variance to guide inventory decisions is common practice, essentially because it leads to efficient strategies for pooling inventories (over products, markets, and periods). Moreover, the variance of aggregated demands is easy to compute, while the shape of the distribution might not be preserved under convolution (e.g., see Eppen and Martin 1988). However, estimating the variance can sometimes be more an art than a science (Fisher and Raman 1996). In addition, when qualitative forecast methods are used, it might be easier to characterize the shape of the demand distribution. It is therefore important to quantify the value of information associated with the variance, in comparison to a more qualitative description of demand.

Figure 1-6 compares the Price of Information under different assumptions on the demand distribution, when $\mu = 100$ and $c = 1$ (without loss of generality). As one adds more restrictions to the shape of the demand distribution, the Price of Information ($PoI$) decreases. Moreover, the $PoI$ generally increases with the profit margin, unless one assumes that the median is equal to the mean, in which case the $PoI$ first decreases to zero (at $r = 1/2$) and increases afterwards. In fact, when the median is known and when $r = 1/2$, it is optimal to order the median, independently of the shape of the demand distribution.

Figure 1-6 also shows that the regret when the demand distribution is known to be unimodal, symmetric, with mean $\mu$, is lower than the regret when the demand distribution is known to have a mean $\mu$ and a coefficient of variation greater than .3. (which is not an uncommon value for practical coefficients of variation). Similarly,

Figure 1-6: Price of Information

the regret when the demand distribution is known to be symmetric with mean $\mu$ is lower than the regret when the demand distribution is known to have a coefficient of variation greater than one. (In fact, the largest coefficient of variation of a nonnegative symmetric demand distribution is 1, obtained with a two-point distribution with equal probability masses at zero and at $2\mu$.)

Therefore, in practice, when it is clear that the demand distribution is symmetric and unimodal, it does not seem worth spending time and efforts for estimating the demand variance. This observation also suggests alternative forecasting techniques: rather than focusing on the estimation of moments of the demand distribution, one could alternatively try to better characterize the shape of the demand distribution.

# 1.6 Extension: Robust Revenue Management

The newsvendor model (1) is at the foundation of stochastic inventory theory. Practical environments are of course more complex than the single-item, single-period model. In this section, we illustrate the potential of the minimax regret criterion for solving practical problems. In particular, we analyze the multi-item newsvendor model with capacity constraints when demand distributions are partially specified.

We consider a single-period stochastic inventory management problem with $M$ items subject to $K$ capacity constraints. There are $\mathbf{b}$ available units of resources. Every sale of item $j$, $j = 1, ..., M$, generates a revenue $p_j$ and consumes $\mathbf{A}_j$ units of resources. For simplicity, we assume that the unit order costs are equal to zero. The demand for item $j$ is unknown but assumed to be distributed according to $D_j$. The problem is to determine the order quantities $\mathbf{y}$ so as to maximize the expected revenues, subject to the capacity constraints, that is

$$
\begin{aligned}
\max_{\mathbf{y}} \quad & \sum_{j=1}^{M} p_j E[\min\{y_j, D_j\}], \\
s.t. \quad & \mathbf{A}\mathbf{y} \leq \mathbf{b}, \\
& \mathbf{y} \geq \mathbf{0}.
\end{aligned}
\tag{1.10}
$$

Problem (1.10) is commonly solved in network revenue management (RM). Consider a network with $K$ resources (e.g., flights, night stays) and $M$ classes of customers characterized by a specific origin-destination and fare (ODF) pair. Customers arrive according to a certain (continuous or discrete) stochastic process over a finite time interval. We seek a policy $P$ that maximizes the expected revenues $\mathbf{p}'E[\mathbf{n}^P]$, where $\mathbf{n}^P$ is the vector of total number of accepted requests when policy $P$ is in use. The policy needs to satisfy (almost surely) the capacity constraints, i.e., $\mathbf{A}\mathbf{n}^P \leq \mathbf{b}$, and the accepted requests are obviously nonnegative and cannot exceed the demand, i.e., $\mathbf{0} \leq \mathbf{n}^P \leq \mathbf{D}$. In addition, the policy is required to be non-anticipating. That is, the acceptance/rejection decision at each time $t$ should be based only on the information acquired up to time $t$. Let $\mathcal{P}$ be the set of non-anticipating policies. The network

RM problem can then be formulated as follows:

$$
\begin{aligned}
\sup_{P \in \mathcal{P}} \quad & \mathbf{p}'E[\mathbf{n}^P], \\
s.t. \quad & \mathbf{A}\mathbf{n}^P \leq \mathbf{b} \quad (w.p.1), \\
& 0 \leq \mathbf{n}^P \leq \mathbf{D} \quad (w.p.1).
\end{aligned}
\tag{1.11}
$$

With a discrete stochastic process of arrivals, problem (1.11) can be formulated as a dynamic program: At each (discrete) time $t$, given the level of available capacity, one needs to decide about accepting or rejecting the requests arriving in period $t$, in order to maximize the total expected revenues until the end of the time horizon. However, the dynamic program is rarely solved in practice because of computational and management limitations. Practical networks are extremely large, leading to the so-called curse of dimensionality. Moreover, the data requirements of the stochastic process of arrivals might be hard to collect. In fact, Boyd and Bilegan (2003) report that "business practitioners conceptualize demand in terms of units of inventory sold and its variability, not in terms of arrival rates and interarrival distributions, and this can lead to confusion."

Therefore, in practice, the network RM problem (1.11) is often solved approximately. In particular, the problem can be simplified by assuming that capacity cannot be reallocated in the future. Under this assumption, one needs to know only the distribution of the total number of requests per class, and not their order of arrivals (i.e., the underlying stochastic process). The capacity is therefore partitioned into $M$ buckets, each corresponding to a different class. For each class $j$, let $y_j$ be the booking limit; the total sales corresponding to class $j$ are thus equal to $\min\{y_j, D_j\}$, i.e., one cannot sell more than the booking limit. The problem is to choose a capacity allocation $\mathbf{y}$ that maximizes the network expected revenue. Under the assumption of partitioned booking limits, problem (1.11) is thus equivalent to (1.10), which is often called the Probabilistic Nonlinear Optimization Problem (PNLP).

However, problem (1.10) presents its own challenges, because it assumes the knowledge of the demand probability distributions. In practice, the demand distribution

for a particular fare class on a particular flight might be difficult to estimate because it is associated with limited sales data and is highly sensitive to changes in the environment (such as price, seasonality, etc.). Moreover, the large size of practical networks[1] motivates the use of simple forecasting techniques. For tractability, it is also traditionally assumed in (1.10) that the demand distributions of different classes of customers are independent and that there are no group arrivals, which is an obvious simplification of reality.

A common way to simplify the PNLP is to replace the stochastic demand $\mathbf{D}$ by its mean $\boldsymbol{\mu}$. In this case, the problem simplifies to a simple deterministic linear optimization problem (DLP), and combines therefore an efficient solution method with modest data requirements:

$$
\begin{aligned}
\max_{\mathbf{y}} \quad & \mathbf{p}'\mathbf{y}, \\
s.t. \quad & \mathbf{A}\mathbf{y} \leq \mathbf{c}, \\
& \mathbf{0} \leq \mathbf{y} \leq \boldsymbol{\mu}.
\end{aligned}
\tag{1.12}
$$

However, the DLP completely ignores the stochastic nature of the demand. In this section, we formulate the robust version of (1.10), based on the knowledge of the range of the demand distribution. The robust formulation of (1.10) is a simple linear optimization problem but captures the stochastic nature of demand. We will then show how the model can be used in an environment where capacity allocations are made dynamically, as a heuristic for solving (1.11).

**Literature review.** The problem of allocating fixed capacity to several classes of customers, each characterized with a certain probability distribution, has been extensively studied in the literature. The multi-item newsvendor problem (1.10), sometimes called the "newsstand problem," was first presented in Hadley and Whitin (1963). The solution of this model is usually based on a Lagrangean relaxation of the capacity constraints (see e.g., Hadley and Whitin 1963, Hodges and Moore

---

[1]Talluri and van Ryzin (2004) report that, for a medium-size airline, about 2 million demand quantities must be forecasted every day.

1970). Even with one capacity constraint, numerical problems might arise when the number of items is large; Nahmias and Schmidt (1984) developed several heuristics to efficiently approximate the optimal order quantities. When the number of capacity constraints is large, advanced nonlinear optimization methods must be used to identify the set of active constraints; e.g., see Lau and Lau (1995, 1996). Vairaktarakis (2000) investigated the robust version of the problem under the maximin and the minimax regret criteria, and solved the problem with a Langrangean relaxation. However, his model did not consider the capacity constraints in the regret maximization.

When classes are not partitioned, i.e., when the same product can be used to satisfy several class demands, the problem of allocating fixed capacity to different classes of customers is usually formulated as a dynamic program (see Talluri and van Ryzin 2004 for a review). The single-resource capacity control problem was introduced by Littlewood (1972) with two classes of customers arriving sequentially, and was subsequently extended to multiple classes of customers, (Brumelle and McGill 1993, Curry 1990, and Wollmer 1992), with general order of arrivals (Robinson 1995). The optimal control is well understood and can be achieved with either nested protection levels, nested booking limits, or bid price tables (see Talluri and van Ryzin 2004 for a review). Several heuristics, such as the EMSR-a (Belobaba 1987) and EMSR-b (Belobaba 1992), also perform well in practice.

The network RM problem is significantly more complex and little is known about the optimal policy. According to Talluri and van Ryzin (2004), network RM creates methodological and operational challenges, both on the forecasting side and on the optimization side. Therefore, in practice, the network RM problem is solved only approximately. Commonly used controls are partitioned booking limits, virtual nesting controls, and bid-price controls (see Williamson 1988, 1992 for a comparison of those approaches). The network RM is often approximated with the DLP (Dror et al. 1988, Bertsimas and Popescu 2003), with possible randomization of demand samples (Talluri and van Ryzin 1999), or with the PNLP (Wollmer 1986). Both problems assume partitioned booking limits. Curry (1990) and Chi (1995) proposed an intermediate approach when fare classes are only nested by itinerary. The controls

obtained with these approximated problems can also be improved using a stochastic gradient algorithm (Bertsimas and de Boer 2005 and van Ryzin and Vulcano 2005).

**Contributions.** In this section, we model the robust network RM problem with partitioned booking limits (1.10), using the maximin (1.1) and the minimax regret (1.2) criteria. In particular, we assume that only the ranges of the demand distribution are specified, that is, demand for class $j$ is known to have its support on the interval $[l_j, u_j]$. Let $\mathcal{D}$ be the set of multivariate probability distributions such that the support of the $j$th marginal distribution lies in the interval $[l_j, u_j]$. In contrast to most RM models, we make no assumption about the independence of classes and implicitly capture group arrivals, sell-ups, buy-downs, and the influence of external factors on the demand distributions. This representation of uncertainty captures the stochastic nature of the problem, while remaining simple to estimate.

All our formulations are polynomial-size linear optimization problems (LP), which makes them attractive for solving practical large networks. Moreover, the minimax regret problem can be used to derive bid prices that capture the stochastic nature of demand at a low computational cost and can be used as controls in a non-partitioned network environment. Our numerical study confirms the potential of robust optimization, and in particular of the minimax regret approach, as a new methodology for solving practical RM problems.

## 1.6.1 Maximin

First, consider the maximin criterion (1.1) applied to the network problem with partitioned booking limits (1.10):

$$\max_{\mathbf{y}:\mathbf{A}\mathbf{y}\leq\mathbf{b},\mathbf{y}\geq 0} \min_{\mathbf{D}\in\mathcal{D}} \mathbf{p}'E[\min\{\mathbf{y},\mathbf{D}\}] \tag{1.13}$$

The revenue is minimized when $D_j = l_j$ w.p.1, for all $j$. The maximin booking limits can be found by solving a problem similar to (1.12), in which the vector of mean demand $\mu$ is replaced with the lower value of the support l.

To reduce the level of conservativeness of the maximin approach, Bertsimas and Sim (2004) proposed to restrict the total amount of variability to a "budget of uncertainty." Indeed, in practice, it is unlikely that all demands differ significantly from their nominal (or expected) values. For each demand $D_j$, let $m_j$ be the nominal value and let $s_j$ the maximum deviation from it. With these notations, the demand has a lower bound $l_j = m_j - s_j$ and an upper bound $u_j = m_j + s_j$. Bertsimas and Sim (2004) suggested to limit the sum of relative deviations from the nominal values to some budget $\Gamma$, that is to require that $\sum_{j=1}^{M} |(d_j - m_j)/s_j| \leq \Gamma$. The parameter $\Gamma$ measures the amount of uncertainty captured by the model. If $\Gamma = 0$, there is no uncertainty and the problem is deterministic (in which all demands $d_j$ are equal to their nominal values $m_j$). At the other extreme, if $\Gamma$ equals the number of classes, there is complete uncertainty on demand, and the worst case occurs when $d_j = m_j - s_j$, $\forall j$. Intermediate values of $\Gamma$ specify a moderate level of conservativeness, while accounting for some uncertainty. For simplicity, we assume that $\Gamma$ is an integer.

The objective is to determine booking limits $\mathbf{y}$ that maximize the worst-case *realized* profit, under the budget constraint, that is

$$\max_{\mathbf{y}:A\mathbf{y}\leq\mathbf{b},\mathbf{y}\geq 0} \min_{\mathbf{d}\in[\mathbf{l},\mathbf{u}]:\sum_{j=1}^{M}|(d_j-m_j)/s_j|\leq\Gamma} \mathbf{p}'\min\{\mathbf{y},\mathbf{d}\} \qquad (1.14)$$

We show next that the maximin network RM problem with a budget of uncertainty can be formulated as a polynomial-size linear optimization problem. For each class $j$, the revenue is minimized only with negative deviations from the nominal values. Therefore, the maximin network RM problem with a budget of uncertainty $\Gamma$ can be simplified to:

$$\max_{\mathbf{y}:A\mathbf{y}\leq\mathbf{b},\mathbf{y}\geq 0} \min_{\boldsymbol{\delta}:0\leq\boldsymbol{\delta}\leq 1,\mathbf{1}'\boldsymbol{\delta}\leq\Gamma} \mathbf{p}'\min\{\mathbf{y},\mathbf{m}-\mathbf{S}\boldsymbol{\delta}\},$$

where $\mathbf{S}$ is a diagonal matrix in which diagonal elements are equal to $s_j$, $j = 1,...,M$.

**Theorem 8.** *The maximin network RM problem with partitioned booking limits, with a budget of uncertainty $\Gamma$ on the sum of relative deviations from the nominal values,*

54

*(1.14), can be formulated as the following linear optimization problem:*

$$\max \quad \alpha\Gamma + \beta'\mathbf{1} + \mathbf{p}'\mathbf{w}$$

$$s.t. \quad \mathbf{Ay} \leq \mathbf{b},$$

$$\mathbf{w} \leq \mathbf{y},$$

$$\mathbf{w} \leq \mathbf{m}, \tag{1.15}$$

$$\mathbf{Wp} + \alpha\mathbf{1} + \beta \leq (\mathbf{M} - \mathbf{S})\mathbf{p},$$

$$\alpha \leq 0, \beta \leq 0, \mathbf{y} \geq \mathbf{0},$$

*where* $\mathbf{M}$, $\mathbf{S}$, *and* $\mathbf{W}$ *are diagonal matrices in which diagonal elements are* $m_j$, $s_j$, *and* $w_j$ *respectively.*

Since the budget of uncertainty is restricted to $\Gamma$, the realized revenue will be no less than the optimal value of (1.14) with at most probability $\exp\left(-\Gamma^2/(2M)\right)$, if the demands distributions are symmetric about their nominal values and independent of each other (Bertsimas and Sim 2004).

## 1.6.2 Minimax Regret

The robust network RM problem with the minimax regret criterion (1.2) can be formulated as follows:

$$\min_{\mathbf{y}:\mathbf{Ay}\leq\mathbf{b},\mathbf{y}\geq\mathbf{0}} \left\{ \max_{\mathbf{D}\in\mathcal{D}} \max_{\mathbf{z}:\mathbf{Az}\leq\mathbf{b},\mathbf{z}\geq\mathbf{0}} E[\rho(\mathbf{z},\mathbf{D},\mathbf{y})] \right\}, \tag{1.16}$$

The regret is defined as the difference in revenues:

$$\rho(\mathbf{z},\mathbf{d},\mathbf{y}) = \mathbf{p}'\min\{\mathbf{z},\mathbf{d}\} - \mathbf{p}'\min\{\mathbf{y},\mathbf{d}\}.$$

Following Section 1.3, we invert the order of the two maximization problems. The next lemma characterizes the worst-case demand distributions and the maximum regret attained with that distribution.

**Lemma 1.** *Fix* $\mathbf{y},\mathbf{z}$.

*(a) The optimal demand distributions are unit impulses.*

*(b) $D_j = \max\{z_j, l_j\}$ w.p.1, for all $j$.*

*(c)* $\max_{d_j \in [l_j, u_j]} \rho_j(z_j, d_j, y_j) = \begin{cases} p_j(z_j - y_j), & \text{if } y_j < l_j, \\ p_j \min\{z_j - l_j, (z_j - y_j)^+\}, & \text{if } y_j \geq l_j. \end{cases}$

From Lemma 1, the maximum regret is thus not necessarily concave in $z_j$. In fact, one can show that the problem of choosing the booking limits $\mathbf{z}$ can be formulated as a mixed-integer linear optimization problem. In order to remain tractable, we introduce the following assumption.

**Assumption 1.** *The booking limits $\mathbf{Z}$ can be randomized and the capacity constraint must hold only in expectation, that is $\mathbf{A}E[\mathbf{Z}] \leq \mathbf{b}$.*

Every booking limit $z_j$ is associated with a demand impulse $D_j = \max\{z_j, l_j\}$ w.p.1 from Lemma 1 (b). Therefore, randomizing over the booking limits $\mathbf{z}$ is equivalent to randomizing over the demand impulses, and then choosing a booking limit after observing the demand *realization*. Under Assumption 1, the minimax regret network RM problem (1.16) is then effectively simplified to

$$\min_{\mathbf{y}: \mathbf{Ay} \leq \mathbf{b}, \mathbf{y} \geq 0} \left\{ \max_{\mathbf{D} \in \mathcal{D}} \max_{\mathbf{Z} \in \mathcal{Z}} E[\rho(\mathbf{Z}, \mathbf{D}, \mathbf{y})] \right\}, \tag{1.17}$$

where $\mathcal{Z}$ is the set of nonnegative multivariate distributions such that $\mathbf{A}E[\mathbf{Z}] \leq \mathbf{b}$. Since $\mathbf{AZ} \leq \mathbf{b}$ w.p.1 $\Rightarrow \mathbf{A}E[\mathbf{Z}] \leq \mathbf{b}$, the regret with this relaxed constraint is an *upper bound* on the maximum regret (1.16).

Since the booking limits $\mathbf{Z}$ are randomized under Assumption 1, the expected regret corresponds to the concave envelope of the function presented in Lemma 1 (c), with at most three breakpoints: at zero, $l_j$, and $u_j$. Therefore, any booking limit $z_j$ can be represented as a convex combination of these three breakpoints. The next theorem shows that the minimax regret network RM (1.17) can be efficiently solved.

**Theorem 9.** *The minimax regret network RM problem with partitioned booking limits, (1.17), where the capacity constraints must hold only in expectation, can be for-*

*mulated as the following linear optimization problem:*

$$\text{min} \qquad \boldsymbol{\pi}'\mathbf{b} + \mathbf{q}'\mathbf{1}$$

$$\text{s.t.} \qquad \mathbf{Ay} \leq \mathbf{b},$$

$$\boldsymbol{\pi}'\mathbf{AU} + \mathbf{q}' \geq \mathbf{p}'(\mathbf{U} - \mathbf{Y}),$$

$$\boldsymbol{\pi}'\mathbf{AL} + \mathbf{q}' \geq \mathbf{0},$$

$$\boldsymbol{\pi}'\mathbf{AL} + \mathbf{q}' \geq \mathbf{p}'(\mathbf{L} - \mathbf{Y}), \qquad (1.18)$$

$$\mathbf{q}' \geq -\mathbf{p}'\mathbf{L},$$

$$\mathbf{q}' \geq -\mathbf{p}'\mathbf{Y},$$

$$\mathbf{y} \leq \mathbf{u},$$

$$\boldsymbol{\pi}, \mathbf{y} \geq \mathbf{0}.$$

*where* $\mathbf{U}, \mathbf{L}, \mathbf{Y}$ *are diagonal matrices in which the diagonal elements correspond to* $\mathbf{u}, \mathbf{l},$ *and* $\mathbf{y}$ *respectively.*

The linear problem presented in Theorem 9 has $(K+2M)$ variables and $(K+6M)$ constraints. In comparison, the DLP has $M$ variables and $K + M$ constraints. The larger size of the problem is the price to pay to capture demand stochasticity.

**Regret with no uncertainty.** There is no uncertainty if $u_j = l_j = \mu_j$ for all $j$; that is, demand is deterministic. In this case, the minimax regret LP (1.18) simplifies to

$$\text{min} \quad -\mathbf{p}'\mathbf{y} + \boldsymbol{\pi}'\mathbf{b} + (\mathbf{p}' - \boldsymbol{\pi}'\mathbf{A})^+\boldsymbol{\mu},$$

$$\text{s.t.} \qquad \mathbf{Ay} \leq \mathbf{b},$$

$$\mathbf{y} \leq \boldsymbol{\mu},$$

$$\boldsymbol{\pi}, \mathbf{y} \geq \mathbf{0}.$$

Because the dual of the Deterministic Linear Optimization Problem (1.12) minimizes $\boldsymbol{\pi}'\mathbf{b} + (\mathbf{p}' - \boldsymbol{\pi}'\mathbf{A})^+\boldsymbol{\mu}$, the above problem's objective is to minimize the duality

gap of the DLP. By strong duality, the maximin regret equals then zero, and the optimal booking limits and shadow prices are the same as those obtained with the DLP.

**Regret under uncertainty.** The range is uninformative if $u_j$ is larger than the largest component of the capacity limit $\mathbf{b}$, and $l_j = 0$. Under these circumstances, there is complete uncertainty about the demand distributions. Since $l_j = 0$, $z_j$ is always greater than $l_j$, and the minimax regret LP (1.18) simplifies to

$$\min \quad \boldsymbol{\pi}'\mathbf{b} + \mathbf{q}'\mathbf{1},$$
$$s.t. \quad \boldsymbol{\pi}'\mathbf{A}\mathbf{U} + \mathbf{q}' \geq \mathbf{p}'(\mathbf{U} - \mathbf{Y}),$$
$$\mathbf{A}\mathbf{y} \leq \mathbf{b},$$
$$\mathbf{y}, \boldsymbol{\pi}, \mathbf{q} \geq \mathbf{0}.$$

The variable $\mathbf{q}$ can be eliminated by making the first constraints tight, that is

$$\min \quad \boldsymbol{\pi}'\mathbf{b} + (\mathbf{p}'(\mathbf{U} - \mathbf{Y}) - \boldsymbol{\pi}'\mathbf{A}\mathbf{U})^+\mathbf{1},$$
$$s.t. \quad \mathbf{A}\mathbf{y} \leq \mathbf{b},$$
$$\mathbf{y}, \boldsymbol{\pi} \geq \mathbf{0}.$$

For a single-leg problem, the optimal booking limits are given in closed form. Because $u_j > b$ by assumption, the worst-case demand is such that for some $j$, $D_j = b$ w.p.1, and $\forall k \neq j$, $D_k = 0$ w.p.1, and has a value $p_j(b - y_j)$. Since there are only $M$ such demand scenarios, they can be enumerated. As a result, the minimax regret LP simplifies to

$$\min \quad \max_j\{p_j(b - y_j)\},$$
$$s.t. \quad \mathbf{1}'\mathbf{y} \leq b,$$
$$\mathbf{y} \geq \mathbf{0}.$$

**Proposition 2.** *In the case of complete uncertainty, the partitioned booking limits that minimize the maximum regret in a single-leg problem are given by:*

$$
y_j = \begin{cases} b - b \frac{t-1}{\sum_{k \le t} \frac{1}{p_k}} \frac{1}{p_j}, & \text{if } j \le t, \\ 0, & \text{if } j > t. \end{cases}
$$

*where $t$ is the largest integer less than $M$ such that $(t-1)/(\sum_{j \le t}(1/p_j)) \le p_t$, or equal to $M$ otherwise.*

Therefore, when booking limits are partitioned, the fraction of capacity allocated to some class is proportional to its fare relative to the harmonic mean of all fares.

For instance, if $b = 10$ and $p_i = 4 - i$ for $i = 1, 2, 3$, then $t = 2$ because $p_2 = 2 \ge (2 - 1)/(1/3 + 1/2) = 6/5$, but $p_3 = 1 < (3 - 1)/(1/3 + 1/2 + 1) = 12/11$. Thus, it is optimal to set $y_1 = 6, y_2 = 4$ and $y_3 = 0$. Notice that these booking limits have been derived without making any assumptions about the demand distribution.

**Regret with knowledge about the moments.** One can also incorporate some moment information into $\mathcal{D}$. However, the minimax regret problem might no longer have a linear formulation, as (1.18), unless the regret function is approximated with a piecewise linear function.

For instance, suppose that, in addition to the range of the distribution $[0, u_j]$, one also knows the mean demand $\mu_j$. Similarly to Lemma 1, it is possible to derive an explicit function for the regret. In particular, using a similar argument as in Theorem 2 in Section 1.4, one can show that the expected regret associated with a pair of booking limits $y_j$ and $z_j$, equals $p_j \min\{1, \mu_j/z_j\}(z_j - y_j)^+$. The function is nonlinear between $\max\{y_j, \mu_j\}$ and $u_j$. Nevertheless, in practice, the regret function $p_j \min\{1, \mu_j/z_j\}(z_j - y_j)^+$ can be approximated with a piecewise linear function, with breakpoints at zero, $\mu_j$, and $u_j$. With this approximation, the minimax regret problem can then be formulated as a linear optimization problem, as in Theorem 9.

**Comparison between the minimax regret approach and the multi-item newsvendor problem with uniform distributions.** With a single item, the

minimax regret problem is the same as the newsvendor problem with a uniform distribution (see Theorem 1). However, with multiple items and capacity constraints, the two approaches differ. Consider for instance a problem with $l_1$, $l_2 > 0$ and a single capacity constraint. If the capacity equals $l_1 + 1$, the solution of (1.10) is such that $y_1 = l_1 + 1$ if $p_1(u_1 - y_1)/(u_1 - l_1) > p_2$, by comparing expected marginal seat revenues. In contrast, the solution of the minimax regret solution has $y_1 = l_1 + 1$ if $p_1 > 2p_2$. Indeed, if the unit is allocated to class 2, the maximum regret equals $p_1 - p_2$ (if $l_1 + 1$ units are demanded for class 1), while if the unit is allocated to class 1, the maximum regret equals $p_2$ (if only $l_1$ units are demanded for class 1).

### 1.6.3   Numerical examples with partitioned classes

**Single-leg problem.**  We first consider the single-leg five-class example proposed by Wollmer (1992). Each class $j$ is characterized by a fare $p_j$, a mean demand $\mu_j$, and a standard deviation $\sigma_j$ (see Table 1.6). In addition, we assume that the only information available about the demand distribution is that it has support in the interval $[l_j, u_j]$, where $l_j$ and $u_j$ correspond to the 5th and 95th percentiles of the gamma distribution with mean $\mu_j$ and standard deviation $\sigma_j$. (Given that $\sigma_j/\mu_j = 1/3$ in this example, we could have taken $l_j = \mu_j - 2\sigma_j$ and $u_j = \mu_j + 2\sigma_j$ considering a normal distribution.) There are 119 available seats available on the aircraft.

Table 1.6: Problem data and partitioned booking limits for the single-leg example of Wollmer (1992)

|   |       | Demand Statistics | | Partitioned Booking Limits | | | | |
|---|-------|---------|---------|--------|---------|------|--------|--------|
| $j$ | $p_j$ | $\mu_j$ | $\sigma_j$ | Regret | Maximin | DLP | PNLP N | PNLP U |
| 1 | 1,050 | 17.326 | 5.775 | 20.1 | 12.2 | 17.3 | 19.5 | 20.8 |
| 2 | 567 | 45.052 | 15.017 | 35.3 | 25.6 | 45.0 | 39.1 | 38.8 |
| 3 | 534 | 39.550 | 13.183 | 29.0 | 23.1 | 39.6 | 32.8 | 32.3 |
| 4 | 520 | 34.018 | 11.339 | 24.2 | 20.5 | 17.1 | 27.6 | 27.0 |
| 5 | 350 | 19.786 | 6.595 | 10.3 | 19.8 | 0.0 | 0.0 | 0.0 |

Table 1.6 displays the partitioned booking limits obtained with the minimax regret LP, the maximin LP with $\Gamma = 5$, the DLP, the PNLP with independent normal

distributions with mean $\mu_j$ and standard deviation $\sigma_j$, and the PNLP with independent uniform distributions over $[l_j, u_j]$. The DLP reserves capacity in priority to the high-fare classes, without exceeding their mean demand. On the other hand, the maximin LP allocates capacity so that the minimum demand $l_j$ is covered for all classes $j$; since $\sum_j l_j < b$ in the example, the remaining capacity is allocated indifferently. In contrast to the maximin and the DLP booking limits, the minimax regret booking limits are similar to those obtained with the PNLPs. However, the minimax regret problem is a simple linear optimization problem, while the PNLPs are non linear.

Tables 1.7 and 1.8 compare the expectations and standard deviations of revenues (computed over 1,000 samples) obtained with the different approaches, under the following demand scenarios: independent gamma distributions with mean $\mu_j$ and standard deviation $\sigma_j$, independent uniform distributions over $[l_j, u_j]$, independent Poisson distributions with mean $\mu_j$, independent two-point distributions with probability 1/4 (resp. 3/4) at $l_j$ and probability 3/4 (resp. 1/4) at $u_j$. Not surprisingly, the maximin booking limits perform poorly since they do not reserve enough capacity for the high-fare demand. In contrast, the minimax regret booking limits are near-optimal under all demand scenarios and generally outperform those obtained with the DLP (both in terms of expected revenues and in terms of standard deviation). The DLP gives larger revenues than the minimax regret approach when the demand distribution is Poisson; in this case, the actual range is $[0, \infty)$ and not $[l_j, u_j]$. Incidentally, the booking limits obtained with both PNLPs are also robust, as they perform consistently well under all demand scenarios.

Table 1.7: Expected revenues for the single-leg example of Wollmer (1992) with partitioned classes

| Demand | Regret | Maximin | DLP | PNLP N | PNLP U |
|---|---|---|---|---|---|
| Gamma | 66,342 | 55,382 | 65,341 | 66,873 | 66,760 |
| Uniform | 67,828 | 56,157 | 67,122 | 68,480 | 68,497 |
| Poisson | 68,870 | 56,474 | 69,194 | 70,478 | 70,190 |
| 2-point $(\frac{1}{4}, \frac{3}{4})$ | 66,205 | 54,618 | 66,190 | 66,657 | 67,125 |
| 2-point $(\frac{3}{4}, \frac{1}{4})$ | 53,346 | 49,345 | 50,792 | 51,139 | 51,292 |

Figure 1-7 illustrates that the regret is first increasing and then decreasing with

Table 1.8: Standard deviations of revenues for the single-leg example of Wollmer (1992) with partitioned classes

| Demand | Regret | Maximin | DLP | PNLP N | PNLP U |
|---|---|---|---|---|---|
| Gamma | 5,061 | 2,202 | 6,254 | 5,638 | 5,840 |
| Uniform | 4,429 | 1,368 | 5,570 | 4,893 | 5,198 |
| Poisson | 3,457 | 1,197 | 3,729 | 3,505 | 3,821 |
| 2-point $(\frac{1}{4}, \frac{3}{4})$ | 2,738 | 2,266 | 2,908 | 2,841 | 2,809 |
| 2-point $(\frac{3}{4}, \frac{1}{4})$ | 2,288 | 2,147 | 2,360 | 2,217 | 2,204 |

capacity. In fact, as capacity becomes large, most of the demand is covered (since it is bounded from above) and the opportunity cost is reduced. The regret does not evolve smoothly: it is sometimes constant and sometimes sensitive to the level of available capacity. This reflects the nonlinear nature of the minimax regret approach.



Figure 1-7: Evolution of the regret as capacity increases

Finally, Table 1.9 shows that both the maximum and average difference in revenues between the minimax regret LP and the PNLP with normal distribution is fairly insensitive to the choice of the values of $l_j$ and $u_j$, although the natural choice of the 5th and 95th percentiles dominates.

**Network problem.** Let us now analyze the five-node hub-and-spoke problem proposed by Williamson (1988). Four cities are connected with a hub. Considering all

Table 1.9: Sensitivity to range

| Percentiles of the distribution | [.01, .99] | [.05, .95] | [.1, .9] | [.2, .8] |
|---|---|---|---|---|
| Maximum difference | 2.37% | 1.61% | 1.77% | 2.81% |
| Average difference | 1.28% | 0.82% | 0.87% | 1.29% |

possible origin-destination pairs, there is a total of 20 itineraries on 8 legs. There is one aircraft per leg, and each aircraft has a capacity of 150 seats. There are four classes per itinerary. Each class on a given itinerary is associated with a mean demand, a standard deviation, and a fare. Williamson considered three demand scenarios, corresponding to low, medium, and high demand. These scenarios differ not only with respect to the mean demand but also with respect to the coefficient of variation. The problem data can be found in Tables 5.1-5.3 in Williamson (1988).

When classes are partitioned, the minimax regret approach performs better than all other distribution-free approaches, as illustrated in Table 1.10 for the high demand scenario. In fact, the minimax regret approach is almost optimal and is more tractable than the PNLPs. As in the single-leg example, the PNLP with normal or uniform demand distributions performs consistently well across the different demand scenarios. Figure 1-8 illustrates that the minimax regret approach performs well independently

Table 1.10: Expected revenues for the high-demand network example of Williamson (1988) with partitioned classes

| Demand | Regret | Maximin | DLP | PNLP N | PNLP U |
|---|---|---|---|---|---|
| Gamma | 147,010 | 133,700 | 144,390 | 148,020 | 147,270 |
| Uniform | 156,190 | 136,380 | 148,910 | 155,830 | 156,470 |
| Poisson | 148,020 | 133,430 | 145,750 | 149,310 | 148,420 |
| 2-point $(\frac{1}{4}, \frac{3}{4})$ | 167,300 | 132,820 | 146,240 | 161,510 | 167,270 |
| 2-point $(\frac{3}{4}, \frac{1}{4})$ | 115,420 | 112,410 | 109,590 | 114,390 | 114,960 |

of the level of capacity, assuming independent gamma demand distributions. Finally, Table 1.11 shows that the good performance of the minimax regret is not dependent on the level of demand.

Figure 1-8: Expected revenues with independent gamma demand distributions for the high-demand network example of Williamson (1988)

Table 1.11: Expected revenues for the network example of Williamson (1988) with partitioned classes

| Demand | Regret | Maximin | DLP | PNLP N | PNLP U |
|--------|--------|---------|-----|--------|--------|
| Low | 91,440 | 84,883 | 80,623 | 91,611 | 91,476 |
| Medium | 124,600 | 116,880 | 121,200 | 125,110 | 124,600 |
| High | 147,010 | 133,700 | 144,390 | 148,020 | 147,270 |

## 1.6.4  Application to non-partitioned network problems

We have assumed so far that classes were partitioned. However, in many RM environments, the same product can be used to satisfy the demand of different types of customers. For instance, in the airline industry, seats in the coach cabin are sold at different fares and/or to passengers having different origins and destinations.

The most common controls when classes are not partitioned, are the booking limits and the bid prices. Booking limits are typically defined over a set of classes of customers. In a single-leg problem, it is in fact optimal to *nest* classes by fare; a request from class $j$ is accepted if the sum of requests from classes $j, ..., M$ that have been accepted so far is less than $y_j$. In other words, $b-y_j$ seats are protected from class $j$. In a network environment however, it is not clear how to order classes associated

with different itineraries. Most existing methods (pro-rated EMSR, Displacement-Adjusted Virtual Nesting) reduce the network problem to a collection of single-leg problems, in which the fares are adjusted to account for the network effects (Talluri and van Ryzin 2004).

Bid prices measure the opportunity cost of capacity: In every period $t$ and for every vector of resources $\mathbf{b}$, a price $\pi_k(\mathbf{b}, t)$ is associated with each resource $k$. A request of class $j$ is accepted at time $t$ if and only if the collected fare $p_j$ exceeds the implicit cost of consuming resources, $\mathbf{A}_j'\pi$. Bid price controls are in general not optimal (Talluri and van Ryzin 1998), yet they are widely used in practice. Usually, the shadow prices of the DLP or the Lagrange multipliers of the PNLP are chosen as bid prices (Simpson 1989, Williamson 1992). Talluri and van Ryzin (1999) also suggested to consider as bid prices the expected dual values of a Randomized Linear Optimization Problem (RLP), similar in nature to the DLP (1.12), but in which $\mu$ is replaced by a random demand sample $\mathbf{d}$.

Because the minimax regret LP (1.18), the maximin LP (1.15), the DLP (1.12), the RLP, and the PNLP (1.10) assume that classes are partitioned, there is no guarantee that the controls derived with these models would perform well in an environment where classes are non-partitioned, at least not asymptotically (Cooper 2002, Talluri and van Ryzin 1998). In fact, the controls obtained with the PNLP, which is the most sophisticated model, have poor performance when classes are not partitioned, in comparison to the controls obtained with the DLP, which is the simplest model (Williamson 1992 and de Boer et al. 2002).

**Bid prices.** The variables $\pi$ correspond to the dual variables associated with the constraint $\mathbf{A}E[\mathbf{Z}] \leq \mathbf{b}$ in (1.17). Therefore, the optimal value of $\pi_k$ measures the additional revenue that could be obtained if, in addition to knowing the demand distributions, we were also provided with an additional unit of capacity $b_k$. Accordingly, the optimal value of these variables can be used as a proxy for the marginal value of capacity.

In contrast to the dual values of the DLP, variables $\pi$ capture the stochastic nature

of the demand. In fact, the minimax regret model (1.18) can be rewritten as follows:

$$\min \quad \boldsymbol{\pi}'(\mathbf{b} - \mathbf{AL}) + \mathbf{q}'\mathbf{1} + \mathbf{p}'(\mathbf{L} - \mathbf{Y})^+\mathbf{1},$$

$$s.t. \quad p_j + q_j\frac{1}{l_j} \geq \boldsymbol{\pi}'\mathbf{A}_j \geq p_j\frac{u_j - y_j - (l_j - y_j)^+}{u_j - l_j} - q_j\frac{1}{u_j - l_j} \quad j = 1, ..., N$$

$$\mathbf{Ay} \leq \mathbf{b},$$

$$\mathbf{y} \leq \mathbf{u},$$

$$\mathbf{y}, \boldsymbol{\pi}, \mathbf{q} \geq \mathbf{0}.$$

The opportunity cost of selling $\mathbf{A}_j$ units of capacity is thus compared to the expected marginal seat revenue of class $j$ if the (partitioned) booking limit associated with class $j$ is equal to $y_j$ and the demand for class $j$ is uniformly distributed between $l_j$ and $u_j$. By comparison, with two fare classes, the optimal bid price $\pi$ is equal to $p_1 P(D_1 > y_1)$ when $y_1$ units of capacity are reserved for class 1.

In contrast, the bid-price obtained with the DLP (1.12) is only compared to the full fare. Indeed, when $l_j = u_j = \mu_j$, the constraint involving $\boldsymbol{\pi}'\mathbf{A}_j$ simplifies to

$$p_j + q_j\frac{1}{\mu_j} \geq \boldsymbol{\pi}'\mathbf{A}_j.$$

Based on the example of Wollmer (1992), Figure 1-9 compares the dual prices obtained with the minimax regret LP, the DLP, the PNLP with gamma demand distributions, the maximin LP, and the RLP with 10 samples of gamma-distributed demand (Talluri and van Ryzin 1999). With the maximin LP, the dual price quickly falls to zero when capacity is sufficiently large to cover the minimum demand $l_j$ for all classes. For the DLP, the dual price decreases stepwise, where the $j$th downward jump occurs when the capacity equals the cumulative mean demand of the $j$ highest-fare classes. In contrast, the Lagrange multiplier obtained with the PNLP decreases continuously. The figure illustrates how the shadow price obtained with the minimax regret approach captures well the stochastic nature of demand, despite the fact that it is the output of a linear optimization problem.

66

Figure 1-9: Evolution of the dual prices as capacity increases

**Implementation of a dynamic bid-price control.** In the following, we evaluate the performance of bid-price controls when classes are not partitioned. Bid prices are computed with the following models: the minimax regret LP, the maximin LP with $\Gamma = M$ (for which the shadow prices equal zero if $b > \sum_j l_j$), the DLP, and the RLP with 10 samples of demand. As before, $l_j$ and $u_j$ correspond to the 5th and 95th percentiles of the gamma distribution with mean $\mu_j$ and standard deviation $\sigma_j$.

Algorithm 1 summarizes the way demand samples are generated and bid prices are used as controls, when the bid prices are derived with the DLP. Similar algorithms can be constructed for the RLP, minimax regret, and maximin. For each itinerary $i$, $i = 1, ..., N$, we assume that there are $M$ fare classes, indexed by $j$. The input data consists of the capacity vector $\mathbf{b}$, the resource utilization matrix $\mathbf{A}$, where the $i$th column $\mathbf{A}_i$ represents the amount of resources needed to satisfy a unit request from itinerary $i$. In addition, each fare class $j$ on each itinerary $i$ is associated with a price $p(i, j)$ and a demand distribution $F_{\Gamma(i,j)}$ with support $[l(i, j), u(i, j)]$ and mean $\mu(i, j)$.

Expected revenues are computed over 100 randomly generated demand scenarios. In our numerical study, we assumed that the demands for the different classes are independent and follow a gamma distribution, and that customers arrive sequentially in increasing order of fare. Within a fare class, the order among customers with

different OD pairs was picked randomly.

Bid-prices must be recomputed frequently during the selling season to represent accurately the opportunity cost of capacity utilization (although frequent reoptimization might not always lead to increased revenues, see Cooper 2002). In our numerical study, bid prices are computed every 100 requests (accepted or not). In Algorithm 1, the variable *NberRequests* counts the number of requests that have been introduced since the last reoptimization of the bid prices. Once this variable is equal to 100, bid prices are recomputed.

The demand forecasts that are used to calculate the bid prices are updated using past sales data. If **d** requests have been introduced by the time the bid prices are recomputed, the mean demand is set to the mean residual life $E[D(i,j)|D(i,j) \geq d(i,j)] - d(i,j)$ $\forall j$ in the DLP and the RLP. Moreover, since fare classes are known to arrive sequentially, the mean demand for the lower-fare classes that have been completely addressed are set to zero. Similarly, in the minimax regret LP and the maximin LP, the lower and upper bounds on the demand distribution are also updated to the 5th and 95th fractiles of the residual life distribution respectively. Updating the models with past requests improved the performance of all models in our numerical study.

Finally, a request from a fare class $j$ is accepted only if there is enough capacity to satisfy the minimum demand for the higher-fare classes, sharing the same itinerary. In other words, a request from fare $j$ and itinerary $i$ is accepted if (i) $p(i,j) \geq \pi'\mathbf{A_i}$ and (ii) $\mathbf{b} > \sum_{k<j} l(i,k)\mathbf{A}_i$.

We first consider the single-leg example of Wollmer (1992). Table 1.12 compares the expected revenues when bid prices are computed with the minimax regret LP, maximin LP, the DLP, and the RLP. The minimax regret bid prices slightly outperform the DLP bid prices since they capture the stochastic nature of demand. While being more computationally-friendly, the minimax regret approach also performs slightly better than the RLP by 1%. The maximin approach performs surprisingly well despite the fact that the budget of uncertainty is set to its maximum. In fact, although the bid price is equal to zero most of the time, not all requests are

---

**Algorithm 1** Dynamic Bid-Price Control Simulation when Bid Prices are obtained with the DLP

---

$Revenue = 0$

$NberRequests = 0$

**for all** $j = 1, ..., M$ **do**

    **for all** $i = 1, ..., N$ **do**

        $x = RAND$

        $d(i, j) = F_{\Gamma(i,j)}^{-1}(x)$

    **end for**

    **while** there exists some itinerary $i$ with $d(i, j) > 0$ **do**

        Pick some itinerary $i$ such that $d(i, j) > 0$

        $d(i, j) = d(i, j) - 1$

        Update the mean residual life $\mu(i, j)$

        $NberRequests = NberRequests + 1$

        **if** $NberRequests = 100$ **then**

            Re-Compute Bid Prices based on DLP with mean vector $\boldsymbol{\mu}$ and capacity vector **b**

            $NberRequests = 0$

        **end if**

        **if** $\boldsymbol{\pi}'\mathbf{A}_i \leq p(i, j)$ and $\mathbf{b} > \sum_{k<j} \mathbf{A}_i l(i, k)$ **then**

            $Revenue = Revenue + p(i, j)$

            $\mathbf{b} = \mathbf{b} - \mathbf{A}_i$

        **end if**

    **end while**

    **for all** $i = 1, ..., N$ **do**

        $\mu(i, j) = 0$

    **end for**

**end for**

---

accepted since some level of capacity is protected to cover the minimum demand of the high-fare classes. Frequent reoptimization of the bid prices does not seem to affect the performance.

We next consider the network example proposed by de Boer et al. (2002). The network consists of four cities in series with six possible OD pairs (flights only go in one direction) and three fare classes per itinerary. The fares, mean, and standard deviation of each ODF class can be found in Tables 8 and 9 of de Boer et al. (2002). Each aircraft has 200 seats. We also considered their cases of inflated variance and smaller fare spreads (Tables 10 and 11 in de Boer et al.). Table 1.13 shows that bid prices obtained with the minimax regret LP performs slightly better than the bid

Table 1.12: Expected revenues with bid-price controls in the single-leg example of Wollmer (1992)

| Reopt. Fqcy | Regret | Maximin | DLP | RLP(10) |
|---|---|---|---|---|
| 1/150 | 68,271 | 65,033 | 68,271 | 67,553 |
| 1/100 | 67,249 | 64,311 | 67,258 | 66,804 |
| 1/50 | 67,475 | 64,566 | 67,437 | 66,987 |
| 1/10 | 68,677 | 65,310 | 67,360 | 68,363 |

prices obtained with the DLP or the RLP. In contrast to the single-leg example (Table 1.12), the maximin approach performs poorly because the protection level only covers the minimum demand of the higher-fare classes on the same itinerary.

Table 1.13: Expected revenues with bid-price controls in the network example of de Boer et al. (2002)

| Model | Regret | Maximin | DLP | RLP(10) |
|---|---|---|---|---|
| base case | 76,614 | 56,245 | 75,965 | 75,141 |
| inflated variance | 75,974 | 59,902 | 74,934 | 74,746 |
| smaller fare spread | 64,376 | 55,570 | 64,430 | 63,513 |

The last example that we consider is the five-node hub-and-spoke network proposed by Williamson (1988), under three demand scenarios (see Tables 5.1-5-3 in Williamson 1988). Table 1.14 displays the performance of the bid prices in this example. The minimax regret bid prices perform better than both the DLP and the RLP bid prices by 1%. The maximin LP perform well when capacity is loose: since the maximin bid price is almost equal to zero, the maximin policy consists in accepting all requests, while the other policies impose unnecessary protection levels. Figure

Table 1.14: Expected revenues with bid-price controls in the network example of Williamson (1988)

| Demand | Regret | Maximin | DLP | RLP(10) |
|---|---|---|---|---|
| High | 161,600 | 113,810 | 159,250 | 159,620 |
| Medium | 135,860 | 134,260 | 135,200 | 135,000 |
| Low | 91,692 | 91,692 | 91,692 | 91,692 |

1-10 illustrates that the minimax regret approach performs well independently of the

level of capacity.



Figure 1-10: Expected revenues for the high-demand network example of Williamson (1988)

To test the robustness of our results, we consider the previous examples with random order of arrivals, among all itineraries and classes. Table 1.15 shows that the bid prices obtained with the minimax regret outperform by 1-4% the bid prices obtained with the alternative approaches, which is even greater than with a sequential order of arrivals.

Table 1.15: Expected revenues with bid-price controls with random order of arrivals

| Problem | Regret | Maximin | DLP | RLP(10) |
|---|---|---|---|---|
| Wollmer (1992) | 71,703 | 69,125 | 70,968 | 68,132 |
| de Boer et al. (2002) base case | 79,020 | 74,483 | 76,608 | 74,729 |
| de Boer et al. (2002) inflated variances | 78,914 | 74,650 | 78,413 | 76,744 |
| de Boer et al. (2002) smaller fare spread | 66,255 | 63,712 | 63,300 | 61,821 |
| Williamson (1988) high | 163,350 | 157,300 | 158,060 | 155,700 |
| Williamson (1988) medium | 136,790 | 137,050 | 133,320 | 131,920 |
| Williamson (1988) low | 92,937 | 92,937 | 92,937 | 92,937 |

# 1.7 Conclusions

In this chapter, we propose a robust approach to the newsvendor model with partial demand information. In particular, we derive order quantities that minimize the newsvendor's maximum regret. The minimax regret objective balances the risks of ordering too little against the risk of ordering too much and is consequently less conservative than the maximin approach.

Most of the derived order quantities are simple functions, which makes them attractive for practical application. The minimax regret approach is related to the entropy maximization approach. As a result, the minimax regret order quantity is close (or equal) to the newsvendor solutions with an exponential demand distribution when only the mean is known, with a uniform distribution when only the range is known, and with a normal distribution when only the mean and variance are known and the coefficient of variation is small.

When the mean and the variance are estimated, assuming a normal demand distribution is justified by (a) the Central Limit Theorem, (b) the maximum entropy, (c) the maximin approach, and (d) the minimax regret approach. However, if the coefficient of variation is larger than 0.3, the normality assumption breaks down as the probability of negative demand becomes non-negligible. In this case, the minimax regret approach seems to be the most appropriate approach as it always recommends to order positive and finite quantities.

We also show that recognizing that the demand distribution has a "regular" shape (i.e., symmetric and unimodal) has more informational value than estimating the variance.

Finally, we illustrate the potential of the minimax regret criterion for solving practical problems. In particular, the minimax regret network revenue management problem can be formulated as a linear optimization problem, which makes the approach attractive for solving practical problems, involving large network sizes and many demand uncertainties. In addition, bid prices based on the minimax regret problem capture the stochastic nature of demand while being computationally inex-

pensive. As we showed in our numerical examples, network RM based on the minimax regret bid prices can outperform the traditional approaches.

On the other hand, the maximin approach generally underperforms the minimax regret. However, in uncertain environments, when the decision-maker is risk-averse, the maximin criterion combined with a budget of uncertainty guarantees a minimum level of revenue with a certain probability, which is called the Value-at-Risk (Bertsimas and Sim 2004, Bertsimas and Thiele 2006, Lancaster 2003).

More generally, the field of robust management opens new avenues for modeling uncertainty differently than with probability distributions. As we have seen in this chapter, this new approach can lead to more efficient solution techniques, while requiring less data or making less stringent assumptions.

# Chapter 2

# The Price of Anarchy in Supply Chains

## 2.1 Introduction

The classic newsvendor model (1) assumes a single decision-maker. However, inventory decisions in supply chains are often the result from complex negotiations among supply partners with different, if not conflicting, objectives. In a decentralized supply chain, each partner acts selfishly, by maximizing his/her own profit. As a result, the level of inventory in a decentralized supply chain might not be equal to what would have been optimal to have if the supply chain had been integrated.

In this chapter, we quantify the efficiency of decentralized supply chains that use price-only contracts. Price-only contracts specify a constant per-unit selling price between a buyer and a seller. These contracts are certainly the simplest and the most common mechanisms for governing transactions in supply chains. However, they do not coordinate the supply chain (Cachon 2003), a manifestation of the "double-marginalization" phenomenon: in a decentralized supply chain with two monopolists, two successive markups occur, causing the final price to be higher and the aggregate profits to be lower than if the firms were vertically integrated (Spengler 1950). When demand is stochastic and the retail price is fixed, double marginalization is reflected through reduced inventory levels.

The limited performance of price-only contracts was first investigated in two-stage supply chains by Lariviere and Porteus (2001) and by Cachon and Lariviere (2001). It was then analyzed in more complex supply chains, such as assembly systems (Wang and Gerchak 2003, Gerchak and Wang 2004, and Tomlin 2003), multi-tier assembly systems (Bernstein and DeCroix 2004), distribution systems when demand is stochastic (Anupindi and Bassok 1999, Cachon 2003, Bernstein and Federgruen 2005) or deterministic (Chen, Federgruen and Zheng 2001, Wang 2001). However, no formal analysis has accurately quantified the loss of efficiency associated with these contracts.

To improve coordination in supply chains, various alternative contracts have been proposed: buy-back, revenue sharing, quantity flexibility, sales rebate, and quantity discount contracts (see Cachon 2003 and Lariviere 1999 for a review). However, these more elaborate contracts are typically more costly to negotiate, more complex to administrate, or might create additional moral hazard problems (e.g., see Krishnan, Kapuscinski, and Butz 2004).

Because of the prevalence of price-only contracts in practice and the additional cost of using more elaborate contracts, it is important to quantify the loss of efficiency associated with price-only contracts. Numerical examples in Cachon (2004) show that the relative efficiency of a two-stage decentralized supply chain could be as low as 70-85% for push configurations and 75-90% for pull configurations. Keser and Paleologo (2004) experimentally observed a relative efficiency of 69% in a two-stage push supply chain, but with different wholesale prices and profit allocations that those predicted by the theoretical models (Lariviere and Porteus 2001). But to the best of our knowledge, there has been no formal analysis to quantify the loss of efficiency in decentralized supply chains.

Quantifying the efficiency of a decentralized system relative to the performance of a centralized system has generated a lot of research interest over the past few years. The "Price of Anarchy"–concept introduced by Koutsoupias and Papadim-itriou (1999) and dubbed by Papadimitriou (2001)–measures the ratio of the per-formance of the centralized system over the worst performance of the decentralized

system (corresponding to the worst Nash equilibrium). The Price of Anarchy has been used as a measure of performance for transportation networks (Roughgarden and Tardos 2000, 2002, Correa, Schulz and Stier-Moses 2004, 2005, Perakis 2005), in network resource allocation games (Johari and Tsitsiklis 2004), and network pricing games (Acemoglu and Ozdaglar 2004).

In contrast, there has been little research on quantifying the efficiency of supply chains. Chen et al. (2000) quantified the inefficiency due to the bullwhip effect by comparing the variance of orders with the variance of demand. Martínez-de-Albéniz and Simchi-Levi (2003) computed the Price of Anarchy of a procurement game with option contracts. There are two main differences between option contracts and price-only contracts. First, option contracts are two-dimensional (there is the reservation price and the execution price); as a result, competition is multidimensional and preserves the diversity of suppliers. Second, there is no moral hazard with option contracts, as supply capacity is assumed to be contractible. Finally, Chan and Simchi-Levi (2005) characterized the loss of efficiency in a two-stage supply chain, where the manufacturer chooses the wholesale price and the retailer chooses the retail price under buyback and inventory sharing contracts.

In this chapter, we quantify the impact of double marginalization in supply chains that use price-only contracts. All our bounds are tight and distribution-free (under the restriction that the demand distribution is nonnegative and has an increasing generalized failure rate). Distribution-free bounds are useful when supply chain design decisions must be made without knowledge about demand distribution.

We focus on price-only contracts, as their inability to coordinate the supply chain lies at the foundation of the large body of research on more elaborate contracts (buyback, quantity flexibility, etc.). We assume that the only decisions are the wholesale price and the inventory/capacity levels at each stage. In particular, we assume that the retail price is fixed and that no efforts can be done to improve forecast accuracy, increase sales, or reduce costs. We also ignore the effect of repeated interaction (Anupindi and Bassok 1999), nonzero reservation profits (Lariviere and Porteus 2001, Bernstein and Marx 2005), or renegotiation to a Pareto-improving situation (Ertogral

and Wu 2001, Cachon 2004), as they diminish the impact of double marginalization.

We characterize the efficiency of different supply chain configurations: push or pull inventory positioning, two or more stages, serial or assembly systems, single or multiple competing suppliers, and single or multiple competing retailers. We also validate our findings with a numerical study, measuring supply chain efficiency for commonly-used demand distributions. Our analysis generates the following insights:

1. The loss of efficiency from double marginalization constitutes a major concern: even in a two-stage supply chain, there might be a loss of efficiency of 42%; with more stages, the supply chain performance deteriorates further.

2. The efficiency of price-only contracts generally drops with the number of intermediaries but rallies when competition is introduced. There are however a few exceptions.

3. Focusing on the impact of double marginalization, and assuming everything else being equal, a pull inventory configuration generally outperforms a push configuration. In particular, significant savings can be realized if the supply chain adopts an assemble-to-order production policy, i.e., assembles components after observing the demand.

The chapter is organized as follows. In Section 2.2, we introduce the model framework and characterize the solution of the integrated supply chain. In Section 2.3, we compute the Price of Anarchy of a serial system, first limited to two stages, then consisting of an arbitrary number of stages, in both pull and push configurations. Section 2.4 is devoted to deriving the Price of Anarchy of assembly systems, in which the components are procured from different suppliers. In Sections 2.5 and 2.6, we quantify the efficiency of competition among suppliers and among retailers respectively. Finally, we summarize and discuss our results in Section 2.7, to generate insights into supply chain design.

## 2.2　Model Framework

**Model Notations.**　Consider a supply chain facing the newsvendor problem (1). We assume that the demand distribution $F(x)$ is strictly increasing and continuous, with density $f(x)$, and has an increasing generalized failure rate (IGFR) (see Lariviere and Porteus 2001, Lariviere 2004, and Paul 2005). Let $h(x) = f(x)/\bar{F}(x)$ be the hazard rate and let $g(x) = xh(x)$ be the generalized failure rate, approximating the percentage decrease in the probability of a stockout from increasing the stocking quantity by 1%. The IGFR assumption is sufficient to guarantee a well-behaved (unimodal) problem for the contract initiator in a decentralized setting.

We also define $\ell(x) = (h(x) \int_0^x \bar{F}(\xi)d\xi)/\bar{F}(x)$, roughly representing the percentage decrease in the probability of stockout from increasing the expected sales by 1%. The quantity $\ell(x)$ is increasing if the distribution is IGFR (Cachon 2004).

**Centralized Supply Chain.**　As a benchmark, we consider the centralized (or integrated) supply chain, as if there were a single decision-maker operating the entire supply chain. The level of inventory $y^c$ (for *centralized*) is chosen to maximize the total supply chain expected profits and therefore solves the newsvendor problem (1). That is, $y^c = \bar{F}(r)$ where $r = c/p$.

**The Price of Anarchy.**　The inventory level in a *decentralized* supply chain, denoted by $y^d$, is in general not equal to $y^c$, as each partner optimizes her own profit locally.

To measure the loss of efficiency, we derive a worst-case bound, computed over all IGFR distributions. We restrict our analysis to the class of IGFR distributions to ensure that the decentralized problem has a well-defined solution. This bound is a proxy for the magnitude of the loss of efficiency, without requiring to estimate the demand distribution. In fact, supply chain design, which is essentially strategic, must often be done without knowing the demand distribution (especially if the same supply chain is used for several generations of products), in contrast to inventory decisions, more tactical. All of our bounds are tight; that is, there exists an IGFR demand

Figure 2-1: Decentralized supply-chain framework

distribution for which the supply chain efficiency is characterized by the worst-case ratio.

**Definition 2.** *The Price of Anarchy PoA is the worst-case ratio of the profit achieved by the centralized supply chain over the profit achieved by the decentralized supply chain, that is,*

$$PoA = \sup_{D \in \mathcal{D}} \frac{-cy^c + pE[\min\{y^c, D\}]}{-cy^d + pE[\min\{y^d, D\}]}, \tag{2.1}$$

*where $\mathcal{D}$ is the set of nonnegative demand distributions that have the IGFR property.*

**Decentralized Game Framework.**   In the following, we analyze the performance of the decentralized supply chain depicted in Figure 2-1. The supply chain is divided into three parts: the procurement stage, the bill of materials of an assembly, and the distribution stage. A square represents a product and a circle represents a decision-maker. A solid arrow symbolizes a "goes-into" relationship while a dashed arrow represents a supply/distribution channel.

In the center of the figure, the bill of materials represents an assembly structure, where $N_C$ components are assembled into one end-product. At the procurement

stage, in the left part of the figure, each component can be procured from any of $N_S$ competing suppliers. For simplicity, we assume the same number of suppliers corresponding to each component. Finally, at the distribution stage, the end-product can be sold to the end-market through any of the $N_R$ retailers. Accordingly, the structure of the supply chain is parameterized by the triplet $(N_C, N_S, N_R)$.

To highlight the impact of double marginalization of supplier interdependence through a bill of material, supplier competition, and retailer competition, we analyze several special cases of the general supply chain depicted in Figure 2-1, corresponding to specific values of parameters $N_C$, $N_S$, and $N_R$. In general, the structure of the game depends on whether the parameter values equal 1 or many (i.e., $n$, where $n > 1$). Table 2.1 summarizes the different supply networks that we analyze in the next sections.

Table 2.1: Supply network structures

| Structure | $N_C$ | $N_S$ | $N_R$ | Section |
|---|---|---|---|---|
| serial system | 1 | 1 | 1 | 2.3 |
| assembly system | $n$ | 1 | 1 | 2.4 |
| competitive procurement system | 1 | $n$ | 1 | 2.5 |
| competitive distribution system | 1 | 1 | $n$ | 2.6 |

For each supply chain structure presented in Table 2.1, we consider two inventory configurations, according to the push-pull classification introduced by Cachon (2004): when the downstream (resp., upstream) partner holds the supply chain inventory, the supply chain is said to be operated in a *push* (resp., *pull*) mode. In other words, a push (resp. pull) configuration corresponds to a Make-to-Stock (resp. Make-to-Order) production policy at the downstream stage. To highlight the impact of double marginalization, we assume zero unit production costs and zero salvage values in both configurations.

Consistently with the literature (Lariviere and Porteus 2001, Cachon and Lariviere 2001), we model the problem as a Stackelberg game where a leader proposes a "take-it-or-leave-it" contract to a follower, and we assume perfect information. The timing of the game for a two-stage supply chain is outlined below. Notice that the order of the last two steps of the game depends on the push-pull configuration of the supply

chain.

1. The leader offers the follower a contract specifying the unit wholesale price $w$.

2. The follower accepts the contract if his expected profit is above his reservation profit, assumed to be zero. Otherwise, there is no transaction between the parties.

3. The manufacturer chooses his/her level of inventory at a unit cost $c$.

*Pull:*

4. Demand $D$ is realized.

5. The retailer orders what is needed to meet demand, at a unit cost $w$. Each unit of satisfied demand generates a revenue $p$. If the manufacturer does not have enough inventory, the excess demand is lost at no cost.

*Push:*

4. The retailer places an order, at a unit cost $w$.

5. Demand $D$ is realized. Each unit of satisfied demand generates a revenue $p$. If the retailer does not have enough inventory, the excess demand is lost at no cost.

## 2.3 Serial Supply Chain

We first characterize the efficiency of a two-stage supply chain, under a push or a pull configuration, constituted of a manufacturer and a retailer. We then extend our results to multistage serial structures.

### 2.3.1 Push Serial Supply Chain

In a push supply chain, the inventory is held at the retailer's site, i.e., the retailer orders before observing the demand. We consider two Stackelberg games, depending on who proposes the contract.

**Manufacturer is the Leader.** When the manufacturer chooses the wholesale price $w$ to maximize her profits, she anticipates the retailer's order quantity. That is, she solves the following bilevel optimization problem:

$$\max_w (w-c)y,$$

$$s.t. \quad y = \arg\max_x pE[\min\{x, D\}] - wx, \quad (I.C.),$$

$$pE[\min\{y, D\}] - wy \geq 0, \qquad (I.R.).$$

Consistently with the standard principal-agent models, the incentive compatibility (I.C.) constraint states that the retailer chooses the order quantity to maximize his expected profit, and the individual rationality (I.R.) constraint ensures that the retailer earns at least his reservation profit (assumed to be zero).

Lariviere and Porteus (2001) showed that if the demand distribution is IGFR with a finite mean, the manufacturer's problem is pseudo-concave. In particular, they showed that the manufacturer's optimal sales quantity lied in the interval $[A, \bar{y}]$ where $A \geq 0$ is the lowest value of the support of the distribution, and $\bar{y} = \inf\{x|g(x) \leq 1\}$ maximizes the manufacturer's revenues. They showed that $\bar{y}$ is finite when the mean is finite.

However, requiring the mean to be finite is unnecessary. In fact, the manufacturer's profit is negative when $w < c$, or equivalently when $y > y^c$; on the other hand, the manufacturer can obtain zero profit by charging $w = p$, inducing the retailer not to order. Therefore, the manufacturer's optimal sales quantity lies in the interval $[A, \min\{\bar{y}, y^c\}]$, which is always finite; moreover, on that interval, the manufacturer's profit function is strictly concave. Therefore, the optimal inventory level in the supply chain is either equal to the lowest value of the support of the demand distribution or uniquely determined by the following first-order optimality condition:

$$\bar{F}(y^d)(1 - g(y^d)) = r, \tag{2.2}$$

where $g(x) = xf(x)/\bar{F}(x)$ is the generalized failure rate.

The next theorem quantifies the Price of Anarchy (*PoA*, see (2.1)) of a two-stage push supply chain.

**Theorem 10.** *In a two-stage push supply chain, when the manufacturer is the initiator of the price-only contract,*

$$PoA = r^{\frac{-1}{1-r}} - r^{-1}.$$

The bound is tight. In particular, the worst-case demand distribution is a Pareto distribution, $\bar{F}(x) = S^k x^{-k}$ for $x \geq S$, for some $S > 0$ and with $k = 1 - r$. Such a demand distribution is appropriate when orders of a total of $S$ units have been pre-committed, and all the remaining uncertainty lies on the excess demand above that level $S$.

**Remark 1.** *The Price of Anarchy measures the largest loss of efficiency over the set of IGFR distributions. With more general distributions, the loss of efficiency can be even larger. For instance, in appendix, we consider the two-point demand distribution, equal to $l$ with probability $1 - r$, and equal to $u > l$ with probability $r$, that is associated with a Price of Anarchy of $2 - l/u$. When $l/u \to 0$, this is greater than the bound derived in Theorem 10 when $r > .5$*

Table 2.2 displays the performance of a two-stage push supply chain when the demand distribution is gamma, i.e., $f(x) = x^{k-1} \exp(-x/\theta)/(\Gamma(k)\theta^k)$. Since supply chain efficiency is independent of the scale of the demand (Lariviere and Porteus 2001), we only vary parameter $k$, by taking $k = .1, 1,$ and 10. When $k = 1$, the distribution is exponential. As $k$ increases, the coefficient of variation $1/\sqrt{k}$, the skewness, and the kurtosis decrease, and the gamma distribution looks more like a normal distribution. Also, when the demand distribution is uniform between 0 and $u$, it turns out that $E[\Pi(y^c, D)]/E[\Pi(y^d, D)] = 4/3$, independently of the profit margins and the value of $u$.

84

Table 2.2: Performance of a two-stage push supply chain

| r | PoA | $E[\Pi(y^c, D)]/E[\Pi(y^d, D)]$ | | | $y^c/y^d$ | | | Stockout Probability | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | $k = .1$ | $k = 1$ | $k = 10$ | $k = .1$ | $k = 1$ | $k = 10$ | $k = .1$ | $k = 1$ | $k = 10$ |
| .2 | 2.48 | 1.36 | 1.41 | 1.34 | 2.64 | 2.57 | 1.84 | 0.27 | 0.53 | 0.85 |
| .4 | 2.10 | 1.36 | 1.37 | 1.30 | 2.60 | 2.28 | 1.67 | 0.45 | 0.67 | 0.90 |
| .6 | 1.92 | 1.36 | 1.35 | 1.27 | 2.60 | 2.14 | 1.56 | 0.64 | 0.79 | 0.93 |
| .8 | 1.80 | 1.36 | 1.34 | 1.24 | 2.59 | 2.06 | 1.48 | 0.82 | 0.90 | 0.97 |

As shown in Table 2.2, the worst-case ratio $PoA$ is roughly 30-80% greater than the performance of the supply chain for "common" demand distributions (exponential, gamma, uniform). However, the qualitative behavior remains the same: the performance of the decentralized supply chain is decreasing with the profit margin $1 - r$ and the coefficient of variation $1/\sqrt{k}$, at least after a certain point, as already noted by Lariviere and Porteus (2001). Table 2.2 also shows that the loss of efficiency is caused by both a reduction in inventory levels and a deterioration in service level. Interestingly, the two effects evolve in opposite directions as the coefficient of variation increases (i.e., $k$ decreases) or as the profit margin increases (i.e., $r$ decreases).

**Retailer is the Leader.** If the retailer proposes the contract, there is no moral hazard and the first-best solution is achieved. Formally, the retailer solves the following bilevel optimization problem:

$$\max_{w,y} pE[\min\{y, D\}] - wy,$$
$$s.t \qquad (w - c)y \geq 0, \qquad (I.R.).$$

It is easy to see that the retailer proposes a price $w = c$. As a result, the retailer's problem reduces to (1). Therefore, full coordination is achieved, i.e., $PoA = 1$.

**Multistage Supply Chain.** Let us consider a push serial supply chain, where every stage offers a take-it-or-leave-it contract to the next downstream stage. As double marginalization occurs at every stage, the performance of the decentralized supply chain deteriorates. To highlight the effect of double marginalization, we assume that

85

the intermediaries in the supply chain incur no additional cost, and by consequence, do not add value to the product. We start with characterizing a bound on the level of inventory in the decentralized supply chain $y^d$.

**Lemma 2.** *With an IGFR demand distribution, the optimal inventory level in an n-stage push serial system, is such that $\bar{F}(y^d)(1 - g(y^d))^{n-1} \geq r$.*

Using Lemma 2, the next theorem quantifies the worst-case loss of efficiency.

**Theorem 11.** *In an n-stage serial push supply chain, when every stage offers a contract to the next downstream stage,*

$$
\begin{aligned}
PoA &= \frac{(1 - r^{\frac{-1}{n-1}})(1 - r^{1 - \frac{1}{1 - r^{\frac{1}{n-1}}}})}{1 - r} \\
&\approx \frac{e^{n-1} - 1}{n - 1} + \frac{1}{2}(1 - r)\frac{e^{n-1} - n}{(n-1)^2} + O\left((1 - r)^2\right).
\end{aligned}
$$

The approximation, obtained from a series expansion around $r = 1$, shows that the number of intermediaries has an exponential influence on the value of the Price of Anarchy.

Figure 2-2 depicts the evolution of Price of Anarchy of a serial push supply chain as the profit margin $1 - r$ increases and as the number of intermediaries $n$ increases. Table 2.3 compares the Price of Anarchy with the profit ratio $E[\Pi(y^c, D)]/E[\Pi(y^d, D)]$ of serial supply chains with three and four echelons, under a gamma demand distribution, for different values of the coefficient of variation $1/\sqrt{k}$.[1] Also, when the demand distribution is uniform between 0 and $u$, one can show that the profit ratio equals 16/7=2.29 when there are three echelons, and 64/15=4.27 when there are four echelons, independently of $u$ and the profit margin.

We make the following observations, based on Figure 2-2 and Table 2.3. First, the loss of efficiency generally increases with the profit margin of the supply chain

---

[1]In this numerical study, the set of possible wholesale prices between any two stages was discretized: in particular, only the prices $w = r + (1 - \zeta/1001)r$, for $\zeta = 1, ..., 1000$, were considered. Therefore, some errors might be due to the discretization.

Table 2.3: Performance of an $n$-stage push supply chain

|   | $n = 3$ | | | | $n = 4$ | | | |
|---|---|---|---|---|---|---|---|---|
| $r$ | $PoA$ | $k = .1$ | $k = 1$ | $k = 10$ | $PoA$ | $k = .1$ | $k = 1$ | $k = 10$ |
| .2 | 4.14 | 2.40 | 2.62 | 1.95 | 7.68 | 4.84 | 5.14 | 3.07 |
| .4 | 3.72 | 2.51 | 2.47 | 1.89 | 7.15 | 4.86 | 4.78 | 2.87 |
| .6 | 3.50 | 2.51 | 2.37 | 1.79 | 6.81 | 4.85 | 4.34 | 2.70 |
| .8 | 3.32 | 2.51 | 2.30 | 1.72 | 6.56 | 4.85 | 4.02 | 2.51 |

and with the coefficient of variation. Second, even when the profit margin is small ($r$ close to 1), full coordination is not achieved; in particular, with only two stages, the worst-case supply chain performance $PoA$ is equal to $e - 1$ when $p = c$. Third, the supply chain becomes less efficient when the number of intermediaries in the supply chain increases, reflecting the impact of double marginalization. More specifically, the Price of Anarchy is nearly doubled each time a new intermediary is added.



Figure 2-2: Price of Anarchy of a serial push supply chain

## 2.3.2  Pull Serial Supply Chain

In a pull supply chain, the inventory is held at the manufacturer's site, i.e., the retailer orders after observing the demand. We consider two Stackelberg games, depending on who proposes the contract.

**Retailer is the Leader.** The retailer proposes the manufacturer a wholesale price so as to maximize her expected profits, anticipating the manufacturer's choice of in-

ventory level. Formally, the retailer solves the following bilevel optimization problem:

$$\max_w (p - w) E[\min\{y, D\}],$$

$$s.t. \quad y = \arg\max_x w E[\min\{x, D\}] - cx, \quad (I.C.),$$

$$w E[\min\{y, D\}] - cy \geq 0, \qquad (I.R.).$$

Cachon and Lariviere (2001) showed that the optimal inventory level $y^d$ is either equal to the lowest value of the support of the demand distribution or uniquely determined by the following first-order optimality conditions:

$$\bar{F}(y^d) = r(1 + \ell(y^d)), \qquad (2.3)$$

where $\ell(x) = (\int_0^x \bar{F}(\xi) d\xi) f(x)/(\bar{F}(x))^2$.

The next theorem characterizes the loss of efficiency from decentralizing operations.

**Theorem 12.** *In a two-stage pull supply chain, when the retailer is the initiator of the price-only contract,*

$$PoA = e - 1.$$

*where $e$ approximately equals 2.71828 and is the base of the natural logarithm.*

Equivalently, the loss of efficiency (computed by $1 - 1/PoA$) is 42%. The bound is tight and is attained with the following piecewise Pareto demand distribution, with a breakpoint at $y^d$,

$$\bar{F}(x) = \begin{cases} x^{-l} S^l, & \text{for } S \leq x \leq y^d, \\ x^{-k} T^k, & \text{for } x \geq y^d, \end{cases}$$

where $S^l = (y^d)^l r(1 + k)$, $T^k = (y^d)^k r(1 + k)$, and both $l$ and $k$ tend to zero, but $l$ tends faster to zero than $k$. Roughly speaking, the demand is either equal to zero with probability $1 - r$, or greater than $y^c$ with probability $r$. Thus, in the worst

case, the product is either a success (sales are greater than $y^c$) or a failure. Given the one-shot nature of the game, such a demand distribution is appropriate for the launch of a new product.

The Price of Anarchy in a pull supply chain is lower than that of a push supply chain (Theorem 10) when $r < 1$, and equal when $r = 1$. In fact, in a pull supply chain, the demand risk is spread over the two parties since the transfer of inventory occurs after the demand realization. In contrast, in a push supply chain, only the retailer is directly exposed to the demand risk, as he orders before observing demand. Since there is more risk sharing in a pull supply chain, the individual objectives are more aligned, and the whole supply chain is more efficient.

Table 2.4 compares the Price of Anarchy with the performance of a two-stage pull supply chain, with a gamma distribution and varying coefficient of variation. Also, when the demand is uniform between 0 and $u$, one can show that the ratio takes on values between 1 (when $r \to 0$) and 4/3 (when $r \to 1$), independently of $u$.

Table 2.4: Performance of a two-stage pull supply chain

|  |  | $E[\Pi(y^c, D)]/E[\Pi(y^d, D)]$ | | | $y^c/y^d$ | | | Stockout Probability | | |
|---|---|---|---|---|---|---|---|---|---|---|
| $r$ | $PoA$ | $k = .1$ | $k = 1$ | $k = 10$ | $k = .1$ | $k = 1$ | $k = 10$ | $k = .1$ | $k = 1$ | $k = 10$ |
| .2 | 1.72 | 1.32 | 1.22 | 1.14 | 2.48 | 2.00 | 1.49 | 0.27 | 0.45 | 0.66 |
| .4 | 1.72 | 1.35 | 1.27 | 1.18 | 2.55 | 2.00 | 1.48 | 0.45 | 0.63 | 0.82 |
| .6 | 1.72 | 1.35 | 1.30 | 1.20 | 2.57 | 2.00 | 1.46 | 0.64 | 0.77 | 0.91 |
| .8 | 1.72 | 1.36 | 1.32 | 1.22 | 2.59 | 2.00 | 1.44 | 0.82 | 0.90 | 0.97 |

The Price of Anarchy is in general 30% higher than the profit ratio with gamma and uniform demand distribution. The performance of a pull decentralized supply chain deteriorates when margins become small (in contrast to a push supply chain) or when the coefficient of variation increases (as in a push supply chain).

In addition, comparing Table 2.4 with Table 2.2 reveals that both the inventory and the service levels are improved under a pull configuration relatively to a push configuration.

**Manufacturer is the Leader.** If the manufacturer proposes the contract, there is no moral hazard. Formally, the manufacturer solves the following bilevel optimization

problem:

$$\max_{w,y} wE[\min\{y, D\}] - cy,$$

$$s.t \quad (p - w)E[\min\{y, D\}] \geq 0, \quad (I.R.).$$

It is easy to see that the manufacturer will propose a wholesale price $w = p$. As a result, the manufacturer's problem is equivalent to (1) and the Price of Anarchy is equal to 1.

**Multistage Supply Chain** Let us consider a pull serial supply chain, where every stage offers a take-it-or-leave-it contract to the next upstream stage. As double marginalization occurs at every stage, the performance of the decentralized supply chain decreases with the number of intermediaries. We start with characterizing the level of inventory in the decentralized supply chain, $y^d$.

**Lemma 3.** *With an IGFR demand distribution, the optimal inventory level in an n-stage pull serial system, is such that $\bar{F}(y^d) \geq r(1 + \ell(y^d))^{n-1}$.*

Using Lemma 3, the next theorem quantifies the worst-case loss of efficiency.

**Theorem 13.** *In an n-stage serial pull supply chain, when every stage offers a contract to the next upstream stage,*

$$PoA = \frac{e^{n-1} - 1}{n - 1}.$$

With two stages, the Price of Anarchy is equal to 1.72; with three stages, it goes up to 3.19; with four stages, it raises to 6.36.

As in the two-stage supply chain, the performance of the pull supply chain dominates that of the push supply chain. In particular, the Price of Anarchy of the pull supply chain corresponds to the Price of Anarchy of the push supply chain with zero margin, i.e., $r = 1$ (compare Theorem 13 with the approximation of the Price of Anarchy in a push supply chain in Theorem 11).

90

Table 2.5 compares the performance of the decentralized supply chain, with three or four stages, when the demand distribution is gamma. As in a two-stage supply chain (see Table 2.4), the performance of the supply chain decreases with the coefficient of variation $1/\sqrt{k}$ and increases with the profit margin $1 - r$. Moreover, comparing Tables 2.3 and 2.5 reveals that a pull serial system is more efficient than a push serial system, consistently with the worst-case analysis (except when $n = 4$, $k = .1$, $r \geq .6$ due to roundoff errors).

Table 2.5: Performance of an $n$-stage pull supply chain

| | $n = 3$ | | | | $n = 4$ | | | |
|---|---|---|---|---|---|---|---|---|
| $r$ | $PoA$ | $k = .1$ | $k = 1$ | $k = 10$ | $PoA$ | $k = .1$ | $k = 1$ | $k = 10$ |
| .2 | 3.19 | 2.21 | 1.77 | 1.54 | 6.36 | 4.27 | 2.89 | 2.08 |
| .4 | 3.19 | 2.35 | 2.01 | 1.65 | 6.36 | 4.33 | 3.52 | 2.17 |
| .6 | 3.19 | 2.39 | 2.04 | 1.67 | 6.36 | 5.05 | 3.62 | 2.17 |
| .8 | 3.19 | 2.35 | 2.17 | 1.68 | 6.36 | 5.01 | 3.67 | 2.11 |

## 2.4    Assembly System

In this section, we characterize the efficiency of a two-level assembly system, depicted in Figure 2-3. We consider a manufacturer that purchases for assembly $n$ distinct components from outside suppliers. Following the framework proposed in Figure 2-1, we assume that there are $N_C = n$ components to be assembled and that each of these is procured from exactly $N_S = 1$ supplier. Notice that, in contrast to Section 2.3, the manufacturer is located downstream of the supply chain; accordingly, $N_R = 1$. Each supplier is specialized in the production of a particular component; therefore, suppliers are not competing on the same markets. However, they have strategic interactions as they are linked through the assembly process. Let $c_i$ be the unit production cost of supplier $i$, $i = 1, ..., n$, and $c_0$ be the unit assembly cost of the manufacturer; as before, $p$ is the unit selling price of the end-product.

Since all components are needed for the production of the end-product, the system capacity corresponds to the minimum level of capacity (or inventory) held by each

Figure 2-3: Assembly structure

supplier. As a result, if one supplier is particularly reluctant to build large inventories, it will affect the performance of the entire system.

## 2.4.1 Push Assembly System

We consider first a push supply chain, i.e., where the manufacturer orders before observing the demand. As in a two-stage supply chain, if the manufacturer specifies the terms of the contracts, full efficiency is achieved. On the other hand, if the suppliers specify the contracts, the decentralized supply chain is less efficient than the centralized supply chain. In particular, suppose that each supplier $i$ proposes the manufacturer a unit transfer price $w_i$ in order to maximize her profit. As the manufacturer makes to stock, he faces the standard newsvendor problem. Anticipating the manufacturer's order size, the suppliers propose transfer prices to jointly solve the following bilevel optimization problem:

$$\max_{w_i} (w_i - c_i)y, \qquad\qquad \forall i = 1, ..., n,$$

$$s.t. \quad y = \arg\max_x pE[\min\{x, D\}] - \sum_{i=1}^n w_i x - c_0 x, \quad (I.C.),$$

$$pE[\min\{y, D\}] - \sum_{i=1}^n w_i y - c_0 y \geq 0, \qquad (I.R.).$$

Gerchak and Wang (2004) proved the existence of a Nash equilibrium in the game and showed that the quantity ordered by the manufacturer to each of his suppliers is either

92

equal to the lowest value of the support of the demand distribution or determined by the following equation:

$$\bar{F}(y^d)(1 - g(y^d)) = \frac{c_j + c_0 + \sum_{i \neq j} w_i}{p}, \qquad \forall j = 1, ..., n.$$

Since the left-hand side is constant for every supplier, it turns out that, at the Nash equilibrium, all suppliers make the same profit, i.e., $w_i - c_i = w_j - c_j$, $\forall i, j$. Therefore, one can show that $w_i = c_i + \frac{1}{n}(p\bar{F}(y_d) - \sum_{j=0}^{n} c_i)$, and that the first-order optimality condition reduces to

$$\bar{F}(y^d)(1 - ng(y^d)) = r, \tag{2.4}$$

where $r = \sum_{i=0}^{n} c_i/p$. The next theorem quantifies the loss of efficiency in the assembly system.

**Theorem 14.** *In a push assembly system with $n$ components, when the suppliers offer the price-only contracts,*

$$PoA = \frac{r^{1 - \frac{n}{1-r}} - 1}{n - 1 + r}.$$

As shown in Figure 2-4, the loss of efficiency of price-only contracts can be dramatic in push assembly systems. The performance of the decentralized supply chain quickly falls as the number of components and the profit margins increase. Comparing Figures 2-2 and 2-4 reveals that the Price of Anarchy is larger in push assembly systems with $n$ components than in a $n$-stage push serial structure. Therefore, aiming at reducing the number of stages in the assembly process (e.g., through the production of kits upstream of the assembly stage) can be valuable to mitigate the effect of double marginalization, everything else being equal.

However, the performance of assembly system is not always so poor. Table 2.6 compares the profit ratios $E[\Pi(y^c, D)]/E[\Pi(y^d, D)]$ of an assembly structure with two

Figure 2-4: Price of Anarchy of a push assembly system

or three components, when the demand has a gamma distribution. Similarly, when the demand distribution is uniform between 0 and $u$, one can show that the profit ratio equals $(n + 1)^2/(2n + 1)$, independently of the profit margin and the upper bound $u$; in particular, the profit ratio equals $9/5 = 1.8$ and $16/7 = 2.29$ when there are two and three components to assemble respectively. Unlike the Pareto distribution, which characterizes the worst-case bound, the gamma, exponential, and uniform distributions are associated with efficiency losses that are comparable to (and in fact even lower than) those in serial structures (see Table 2.3). For these distributions, only the number of intermediaries in the supply chain seems to affect its performance, almost independently of the network configuration.

Table 2.6: Performance of an $n$-component push assembly system

| $r$ | $n = 2$ | | | | $n = 3$ | | | |
|-----|---------|---------|--------|--------|---------|---------|--------|--------|
|     | $PoA$   | $k = .1$ | $k = 1$ | $k = 10$ | $PoA$   | $k = .1$ | $k = 1$ | $k = 10$ |
| .2  | 8.48    | 2.35    | 2.09   | 1.53   | 37.54   | 4.26    | 2.81   | 1.66   |
| .4  | 5.34    | 2.32    | 1.96   | 1.47   | 15.86   | 4.18    | 2.58   | 1.59   |
| .6  | 4.20    | 2.33    | 1.89   | 1.43   | 10.26   | 4.17    | 2.45   | 1.53   |
| .8  | 3.58    | 2.32    | 1.84   | 1.38   | 7.76    | 4.18    | 2.35   | 1.47   |

## 2.4.2 Pull Assembly System

Let us consider now a pull assembly system, i.e., where the manufacturer assembles to order but the suppliers decide their inventory level without knowing the exact demand (e.g., because of long lead-times). Suppose that the manufacturer has full negotiating power and offers each of her suppliers a price-only contract. Given he is offered a wholesale price $w_i$, supplier $i$ will hold a certain level of inventory $y_i$. Consequently, the manufacturer will be able to assemble at most $y = \min_{i=1,...,n} y_i$ products. Anticipating the reaction of her suppliers, the manufacturer designs the contracts in order to maximize her profits, that is:

$$\max_{w_1,...,w_n} (p - \sum_{i=1}^n w_i - c_0) E[\min\{\min_i y_i, D\}],$$

$$s.t. \qquad y_i = \arg\max_x w_i E[\min\{x, D, y_{-i}\}] - c_i x, \qquad \forall i = 1,...,n, \quad (I.C.),$$

$$w_i E[\min\{y_i, D, y_{-i}\}] - c_i y_i \geq 0, \qquad \forall i = 1,...,n, \quad (I.R.).$$

The (I.C.) constraint represents the newsvendor problem faced by each supplier, selling the minimum between the demand $D$, his capacity level $y_i$, and the others' capacity levels, denoted by $y_{-i}$. As all components are required to assemble one unit of end-product, the system capacity is the minimum of all supplier's inventory levels. Wang and Gerchak (2003) proved the existence and uniqueness of the Nash equilibrium in the game. They also showed that the optimal level of capacity is either equal to the lowest boundary of the support of the distribution or uniquely determined by the following first-order optimality condition:

$$\bar{F}(y^d) = \frac{\sum_{i=1}^n c_i}{p - c_0} (1 + \ell(y^d)).$$

The quantity defined by (2.5) is the same as the quantity determined by (2.3), as if the supply chain had only one supplier. As a result, the Price of Anarchy is unaffected by the number of suppliers in a pull assembly system.

**Corollary 2.** *In a pull assembly system, when the manufacturer offers the price-only*

*contracts,*

$$PoA = e - 1.$$

Therefore, the performance of a pull assembly system is equal to that of a two-stage pull supply chain. Similarly, an $n$-level pull assembly system is equivalent to an $n$-stage pull serial system (see Bernstein and DeCroix 2004). While the performance of a push assembly system dramatically falls as the number of components increases (see Figure 2-4), the Price of Anarchy of a pull assembly system remains equal to 1.72, independently of the profit margin and of the number of components. Therefore, when the profit margins are large and the bill of materials is complex, it is highly profitable, from a pure double-marginalization perspective, to operate the assembly in a pull instead of a push mode. Part of Dell's success can be explained by its Assemble-to-Order production policy (Dell and Magretta 1998); in addition to lowering system inventories and making the supply chain more responsive to demand variations, this configuration also mitigates the efficiency losses due to double marginalization.

## 2.5 Competitive Procurement System

In this section, we quantify the impact of competition among suppliers on supply chain efficiency. In general, competition improves system efficiency, but the degree of improvement depends on the supply chain structure and the level of competition.

We consider a competitive procurement system similar to the one displayed in Figure 2-5. Following the framework introduced in Figure 2-1, we assume that the manufacturer needs only $N_C = 1$ component, purchased from any of $N_S$ suppliers, to make a unit of final product. All suppliers are identical, i.e., procure the same component and have the same unit cost $c$. As before, the manufacturer is downstream of the supply chain and sells the final product at a unit price $p$; accordingly, $N_R = 1$.

96

Figure 2-5: Competitive procurement system

## 2.5.1 Push Procurement System

In a push supply chain, the manufacturer bears all risk of the supply chain (i.e., is a newsvendor). Let us assume that the suppliers choose the wholesale prices.

Facing different wholesale prices $w_i$'s for the same homogeneous component, the manufacturer orders from the cheapest supplier. Since the suppliers bear no risk, the cheapest supplier will not ration her capacity and will deliver the exact amount ordered. If several suppliers offer the same low price, we assume that the manufacturer splits his order equally among them. This game is exactly a Bertrand competition (Bertrand 1883). Accordingly, there exists a unique Nash equilibrium in which all suppliers offer the same wholesale price $c$.

As a result, there is no loss of efficiency in this supply chain, i.e., $PoA = 1$. It is interesting to notice that the push configuration in a serial supply chain leads to the poorest performance, but that full efficiency is achieved as soon as competition is introduced. In fact, one might think that the opposite situation would happen: as the number of suppliers (leaders) increases, the manufacturer (follower) has less and less surplus, and supply chain efficiency decreases. However, since the leaders compete against each other, they sacrifice the surplus they extract from the manufacturer to remain competitive, and full coordination is achieved.

Similarly, when the manufacturer is the leader, she offers a wholesale price $c$ to each of her suppliers, and the supply chain operates at full efficiency.

## 2.5.2 Pull Procurement System

In a pull supply chain, the suppliers bear the demand risk. The manufacturer is in general not able to procure everything she needs from the same supplier, as this supplier might hold limited inventory. In fact, the cheapest supplier has the least incentives to hold large inventories, as his margins are the thinnest. By increasing the number of suppliers, the manufacturer has more sourcing options and risks less to be held up by one of her suppliers.

Let us assume that the manufacturer initiates the contracts with her suppliers. Anticipating their limited incentives for holding large inventories, the manufacturer offers different wholesale prices to her suppliers. After demand is realized, the manufacturer will procure what she needs from her suppliers in increasing order of wholesale prices. Suppose without loss of generality that supplier $i$ is proposed the $i$th smallest wholesale price and let us denote by $\bar{y}_i$ the cumulative inventory from supplier 1 to supplier $i$, i.e., $\bar{y}_i = \sum_{j \leq i} y_i$. Therefore, the manufacturer will order from supplier $i$ only if demand $D$ exceeds the cumulative inventory level of suppliers 1 to $i - 1$, i.e., if $D \geq \bar{y}_{i-1}$.

Thus, the manufacturer offers different wholesale prices $w_i$'s to maximize her profits, that is:

$$\max_{w_1,...,w_n} \sum_{i=1}^{n} (p - w_i)(E[\min\{D, \bar{y}_i\}|\bar{y}_{i-1} \leq D] - \bar{y}_{i-1})P(\bar{y}_{i-1} \leq D),$$

$$s.t. \qquad \bar{y}_i = \arg\max_{x \geq \bar{y}_{i-1}} w_i(E[\min\{D, x\}|\bar{y}_{i-1} \leq D] - \bar{y}_{i-1})P(\bar{y}_{i-1} \leq D)$$

$$-c(x - \bar{y}_{i-1}), \qquad\qquad \forall i = 1, ..., n, \ (I.C),$$

$$w_i(E[\min\{D, \bar{y}_i\}|\bar{y}_{i-1} \leq D] - \bar{y}_{i-1})P(\bar{y}_{i-1} \leq D) - c(\bar{y}_i - \bar{y}_{i-1}) \geq 0,$$

$$\forall i = 1, ..., n, \ (I.R.),$$

$$w_i \leq w_{i+1}, \qquad\qquad \forall i = 1, ..., n - 1,$$

with $\bar{y}_0 = 0$. The (I.R.) constraints are always satisfied if $\bar{y}_i \geq \bar{y}_{i-1}$ for all $i$. The (I.C.) constraints represent the newsvendor problem faced by each supplier, conditional on the fact that demand exceeds $\bar{y}_{i-1}$ for supplier $i$ to receive an order. The optimal

solution to (I.C.) is such that the cumulative inventory for suppliers $1, ..., i$ equals $\bar{y}_i = \bar{F}^{-1}(c/w_i)$. Because there is a one-to-one correspondence between wholesale prices $w_i$'s and cumulative inventory levels $\bar{y}_i$'s, the manufacturer can equivalently choose the induced inventory levels $\bar{y}_i$'s. Plugging $w_i = c/\bar{F}(\bar{y}_i)$ into the objective function and integrating by parts simplifies the manufacturer's problem to:

$$\max_{\bar{y}_1,...,\bar{y}_n} \sum_{i=1}^{n} (p - \frac{c}{\bar{F}(\bar{y}_i)}) \int_{\bar{y}_{i-1}}^{\bar{y}_i} \bar{F}(x)dx,$$

where $\bar{y}_0 = 0$.

The optimal cumulative inventory levels in the competitive procurement system can then be determined from the following set of first-order optimality conditions:

$$\bar{F}(\bar{y}_1^d) = \bar{F}(\bar{y}_2^d) \left\{ 1 + \frac{f(\bar{y}_1^d)}{(\bar{F}(\bar{y}_1^d))^2} \int_0^{\bar{y}_1^d} \bar{F}(x)dx \right\},$$

$$\bar{F}(\bar{y}_i^d) = \bar{F}(\bar{y}_{i+1}^d) \left\{ 1 + \frac{f(\bar{y}_i^d)}{(\bar{F}(\bar{y}_i^d))^2} \int_{\bar{y}_{i-1}^d}^{\bar{y}_i^d} \bar{F}(x)dx \right\}, \ \forall i = 2, ..., n-1,$$

$$\bar{F}(\bar{y}_n^d) = \frac{c}{p} \left\{ 1 + \frac{f(\bar{y}_n^d)}{(\bar{F}(\bar{y}_n^d))^2} \int_{\bar{y}_{n-1}^d}^{\bar{y}_n^d} \bar{F}(x)dx \right\}.$$

As we show in Lemma 7 in the appendix, the first-order conditions are also sufficient for IGFR demand distributions with increasing $f'(x)\bar{F}(x)/(f(x))^2$; for these distributions, the objective function is concave. In particular, the exponential, the Pareto and the uniform satisfy this condition, so as the Weibull when the shape parameter is less than one. On the other hand, the normal distribution does not satisfy this condition, and there are some instances where the objective is not concave (e.g., when the demand distribution is $N(10, 2)$, with $n = 2$, $p = 2$, $c = 1$, the profit function is only pseudoconcave).

The next theorem quantifies the loss of efficiency in a procurement system. Adding competition improves the performance of the decentralized supply chain. However, full efficiency is only achieved in the limit, when the number of competing suppliers becomes large.

**Theorem 15.** *In a pull procurement system, when the manufacturer offers price-only*

*contracts to n competitive suppliers, then*

$$PoA = \frac{e^{\mathcal{E}^{n-1}(1)} - 1}{\mathcal{E}^{n-1}(1)},$$

*where* $\mathcal{E}(x) = 1 - \exp(-x)$ *and* $\mathcal{E}^n(x)$ *is the n-th composition of the function with itself, i.e.,* $\mathcal{E}^n(x) = \mathcal{E}(\mathcal{E}^{n-1}(x))$ *for* $n > 1$ *and* $\mathcal{E}^0(x) = 1$.

With one supplier, the Price of Anarchy equals 1.72; with two suppliers, the Price of Anarchy goes down to 1.39; with three suppliers, it lowers to 1.28. In the limit, with an infinite number of competing suppliers, the decentralized supply chain is as efficient as the integrated supply chain.

Table 2.7 depicts the performance of competitive procurement systems with two or three suppliers, when the demand distribution is gamma. Comparing Table 2.7 with Table 2.4 illustrates how competition gradually improves supply chain efficiency, consistently with the worst-case analysis.

Table 2.7: Performance of an $n$-supplier competitive procurement system

| $r$ | $n = 2$ | | | | $n = 3$ | | | |
|---|---|---|---|---|---|---|---|---|
| | $PoA$ | $k = .1$ | $k = 1$ | $k = 10$ | $PoA$ | $k = .1$ | $k = 1$ | $k = 10$ |
| .2 | 1.39 | 1.13 | 1.08 | 1.04 | 1.28 | 1.07 | 1.04 | 1.02 |
| .4 | 1.39 | 1.14 | 1.10 | 1.06 | 1.28 | 1.08 | 1.05 | 1.03 |
| .6 | 1.39 | 1.12 | 1.11 | 1.06 | 1.28 | 1.10 | 1.06 | 1.03 |
| .8 | 1.39 | 1.20 | 1.12 | 1.07 | 1.28 | 1.20 | 1.06 | 1.03 |

On the other hand, when the suppliers are the leaders and offer the manufacturer wholesale-price contracts, there exists no pure-strategy Nash equilibrium (see Bryant 1980 for an analogous game). The manufacturer always orders in priority from the cheapest supplier; if two suppliers offer the same price, we assume that the manufacturer orders an equal amount from each of them. Consider two competitors $i$ and $j$ with prices $w_i < w_j$. By offering a price $w_j - \epsilon$ and holding the same level of inventory as supplier $j$, supplier $i$ will make more profit than supplier $j$ with positive probability, since the manufacturer always gives priority to the cheapest supplier. If supplier $i$ offers $w_j - \epsilon$, supplier $j$ can increase its profits by lowering its price to $w_j - 2\epsilon$

and holding the same level of inventory as supplier $i$, contradicting that $w_i < w_j$. Similarly, if $w_i = w_j > c$, supplier $i$ can decrease her price to $w_j - \epsilon$ and make more profit. If $w_i = w_j = c$, supplier $i$ can increase her price and make positive profit.

However, there exists a mixed-strategy Nash equilibrium in that game if the demand distribution is continuous with a finite mean (Dasgupta and Maskin 1986). For each supplier, the feasible set of wholesale prices is the closed interval $[c, p]$. The payoff of supplier $i$ is continuous, except on the set of points for which $w_i = w_j$ for some $j \neq i$. It is bounded since the demand has a finite mean. At the points of discontinuity $w_i = w_j > c$, supplier $i$ can increase its profit by lowering its price; hence, its payoff is everywhere lower semicontinuous. Finally, the sum of profits of all suppliers is continuous since discontinuous changes of orders from one supplier to the other occurs only when both suppliers earn the same profit.

## 2.6    Competitive Distribution System

In this section, we quantify the impact of competition among retailers on supply chain efficiency. In particular, we consider a distribution system, in which $N_R = n > 1$ identical retailers compete for the same demand, each with a unit selling price $p$, as shown in Figure 2-6. In order to compare the performance of the competitive system with an integrated supply chain (1), we assume that the aggregate demand has a distribution function $F$ and is allocated to the retailers according to some rule, as in Lippman and McCardle (1997). All retailers are served by the same manufacturer, producing the same product at a unit cost $c$; accordingly, $N_S = N_C = 1$. Notice that, in contrast to Sections 2.4 and 2.5, the manufacturer is located upstream in the supply chain.

### 2.6.1    Push Distribution System

We first analyze the case where the inventory is located downstream, i.e., at the retailers' sites, assuming that the manufacturer proposes the contracts. We consider two different allocation rules: herd behavior and splitting proportional to inventories.

Manufacturer

Retailers

Figure 2-6: Competitive distribution system

**Herd behavior.** Suppose that customers move as a herd and visit one retailer at a time, in a certain order, until the total demand is satisfied. The order between retailers $\pi = (\pi_1, ..., \pi_n)$ is chosen randomly, and all permutations are equally likely. (See Lippman and McCardle 1997.)

The manufacturer proposes a set of prices $w_1, ..., w_n$ to maximize her expected profits, given that the retailers face the demand risk associated with the customers' herd behavior. Formally, it solves the following problem:

$$\max_{w_1,...,w_n} \sum_{i=1}^{n} (w_i - c) y_i,$$

$$s.t. \quad y_i = \arg\max_x p \sum_{k=1}^{n} \sum_{\pi:i=\pi_k} \frac{1}{n!} E[\min\{x, (D - \sum_{j>i} y_{\pi_j})^+\}] - w_i x, \quad \forall i \ (I.C.),$$

$$p \sum_{k=1}^{n} \sum_{\pi:i=\pi_k} \frac{1}{n!} E[\min\{y_i, (D - \sum_{j>k} y_{\pi_j})^+\}] - w_i y_i \geq 0, \qquad \forall i, \ (I.R.),$$

where the retailers choose their order quantity to maximize their expected profits, taken over all possible permutations.

The retailers are symmetric because all permutations are equally likely. Therefore, if they are offered the same wholesale prices, they will order the same quantity. Given the symmetry, the manufacturer offers the same price to all retailers as they all equally contribute to her profits. Let $y$ be the order quantity of each retailer. The retailer's problem then simplifies to maximize $p \sum_{k=1}^{n} \frac{1}{n} E[\min\{x, (D - ky)^+\}] - wx$, and the optimal order quantity is uniquely determined by $(1/n) \sum_{k=1}^{n} \bar{F}(ky) = w/p$ (Lippman

102

and McCardle 1997). Replacing the wholesale price with the induced order quantity in the manufacturer's profit function, it turns out that the optimal order quantity of each retailer is either equal to the lowest value of the support of the demand distribution or uniquely determined by

$$\sum_{k=1}^{n} \bar{F}(ky)(1 - g(ky)) = r \tag{2.5}$$

The next theorem characterizes the loss of efficiency in this competitive distribution system.

**Theorem 16.** *In a push distribution system with $n$ retailers and herd behavior, when the manufacturer offers the price-only contracts, then*

$$
\begin{aligned}
PoA &= \frac{n}{1 - n(1 - r)}(r^{\frac{n(1-r)-1}{n(1-r)}} - 1) \\
&\approx n(e^{1/n} - 1) + (1 - r)(n^2(e^{1/n} - 1) - e^{1/n}(n - \frac{1}{2})) + O\left((1 - r)^2\right).
\end{aligned}
$$

The approximation, obtained from expanding the Taylor series of $PoA$ around $r = 1$, shows the magnitude of improvement when the number of retailers increases. Figure 2-7 depicts how supply chain efficiency increases with the number of retailers.



Figure 2-7: Price of Anarchy of a push competitive distribution system

103

**Splitting Proportional to Inventories.** Suppose now that the aggregate demand is allocated to the retailers proportionally to their inventory levels (Cachon 2003). There is no reallocation of the unmet demand. Specifically, if retailer $i$ holds $y_i$ units and there are $y$ units in the system, i.e., $y = \sum_j y_j$, the fraction of demand observed by retailer $i$ is $y_i/y$. Consequently, the manufacturer solves the following optimization problem:

$$\max_w \sum_{i=1}^n (w - c)y_i,$$

$$s.t. \quad y_i = \arg\max_x pE[\min\{q, D\frac{x}{x+\sum_{j\neq i} y_j}\}] - w_i x, \quad \forall i = 1, ..., n, \quad (I.C.),$$

$$pE[\min\{y_i, D\frac{y_i}{\sum_j y_j}\}] - w_i y_i \geq 0, \quad \forall i = 1, ..., n, \quad (I.R.).$$

For a given wholesale price $w$, Cachon (2003) showed the existence and uniqueness of the Nash equilibrium of the subgame. He also showed that the optimal order level $y^d$ is determined by the following first-order optimality condition:

$$\bar{F}(y^d)(1 - \frac{1}{n}g(y^d)) = r. \tag{2.6}$$

The next theorem quantifies the loss of efficiency of such a distribution system.

**Theorem 17.** *In a push distribution system with $n$ retailers, when the manufacturer offers the price-only contracts and that aggregate demand is allocated proportionally to the inventory levels, then*

$$PoA = \frac{n}{1 - n(1 - r)}(r^{\frac{n(1-r)-1}{n(1-r)}} - 1).$$

The Price of Anarchy is the same as if customers adopt a herd behavior (see Theorem 16), although the optimality conditions (2.5) and (2.6) are different. Therefore, this bound is robust in the sense that it does not seem to depend on the allocation/reallocation mechanism (as long as there is competition among retailers).

Table 2.8 compares the Price of Anarchy to the performance of distribution sys-

tems with splitting proportional to inventory levels when demand has a gamma distribution.

Table 2.8: Performance of an $n$-retailer competitive distribution system with splitting proportional to inventories

| $r$ | $n = 2$ | | | | $n = 3$ | | | |
|---|---|---|---|---|---|---|---|---|
| | $PoA$ | $k = .1$ | $k = 1$ | $k = 10$ | $PoA$ | $k = .1$ | $k = 1$ | $k = 10$ |
| .2 | 1.51 | 1.10 | 1.13 | 1.17 | 1.30 | 1.05 | 1.06 | 1.10 |
| .4 | 1.42 | 1.10 | 1.13 | 1.16 | 1.25 | 1.05 | 1.06 | 1.10 |
| .6 | 1.36 | 1.10 | 1.13 | 1.14 | 1.22 | 1.05 | 1.06 | 1.09 |
| .8 | 1.32 | 1.10 | 1.12 | 1.14 | 1.20 | 1.05 | 1.06 | 1.09 |

## 2.6.2 Pull Distribution System

In a pull supply chain, the manufacturer bears all demand risk. We consider two allocation schemes: herd behavior and uniform split. (The allocation rule based on the relative inventory position is irrelevant here since retailers order after observing the demand.)

**Herd behavior.** Suppose that all customers visit the same retailer. Since the retailer can order anything she wants from the manufacturer, within the limits of the manufacturer's inventory, she will sell the minimum between the demand and the manufacturer's inventory. In contrast to a push distribution system with herd behavior, customers will not visit more than one retailer. Suppose in addition that each retailer has an equal chance of facing the total demand. Then, retailers choose the wholesale prices they offer to the manufacturer, anticipating the manufacturer's inventory decision, by solving the following problem:

$$\max_{w_i} (1/n)(p - w_i)E[\min\{y, D\}], \qquad \forall i = 1, ..., n,$$

$$s.t. \quad y = \arg\max_x \sum_{i=1}^{n} w_i/nE[\min\{x, D\}] - cx, \quad (I.C.),$$

$$\sum_{i=1}^{n} w_i/nE[\min\{y, D\}] - cy \geq 0, \qquad (I.R.).$$

Therefore, the manufacturer faces a newsvendor problem, with a selling price equal to the average of the proposed wholesale prices. Accordingly, he chooses an inventory level such that the stockout probability $\bar{F}(y)$ equals $cn/\sum_i w_i$.

Given the wholesale prices of the other retailers, each retailer $i$ chooses a wholesale price $w_i$, or equivalently, the induced order quantity $y$ such that $\bar{F}(y)(1+\sum_{j\neq i} w_j/p) = rn(1 + \ell(y))$. If the demand distribution is IGFR, each retailer's profit function is concave. Moreover, the set of possible wholesale prices, $[c, p]$ is compact and convex. As a result, there exists a pure-strategy symmetric Nash equilibrium. Since the right-hand side is the same for all retailers, it turns out that all retailers choose the same wholesale price. Replacing each $w_j$ with $c/\bar{F}(y)$ in the optimality condition simplifies to

$$\bar{F}(y^d) = r(1 + n\ell(y^d)). \tag{2.7}$$

The next theorem characterizes the efficiency of this pull distribution system.

**Theorem 18.** *In a pull distribution system with $n$ retailers with herd behavior, when the retailers propose the price-only contracts, then*

$$PoA = \frac{1}{n}(e^n - 1).$$

With one retailer, $PoA = 1.72$ as in Theorem 12. In contrast to the other competitive systems analyzed so far (competitive procurement, push distribution), competition increases inefficiency in this system. With two retailers, $PoA = 3.19$; with three retailers, $PoA = 6.36$. In fact, because the manufacturer bases his inventory decision on the average wholesale price, retailers have less incentives to discount their wholesale prices, as their impact on the inventory decision is diluted.

**Uniform Allocation.** Suppose now that the aggregate demand is split uniformly among retailers. That is, if the aggregate demand is $D$, each retailer observes a

106

demand $D/n$. Let us assume that the retailers are the leaders and propose the wholesale prices to the manufacturer. After observing the demand, the retailers place their orders to the manufacturer. If the wholesale prices are different, the manufacturer serves the retailers in decreasing order of wholesale prices. In case of equal wholesale prices, the manufacturer allocates an equal amount to each retailer.

If there are no potential entrants, there exists no pure-strategy Nash equilibrium. For instance, consider two competitors $i$ and $j$ with prices $w_j < w_i$. By offering a price $w_j + \epsilon$ and ordering the same amount as retailer $j$, retailer $i$ will make more profit than retailer $j$ with positive probability, since the manufacturer always gives priority to the retailer with the highest price. If retailer $i$ offers $w_j + \epsilon$, retailer $j$ will increase its price to $w_j + 2\epsilon$, contradicting that $w_j < w_i$. Similarly, if $w_i = w_j < p$, retailer $i$ can increase her price to $w_j + \epsilon$ and make more profit. If $w_i = w_j = p$, retailer $i$ can decrease her price and make positive profit. Therefore, there exists no pure-strategy Nash equilibrium. However, there exists a mixed-strategy Nash equilibrium (see Dasgupta and Maskin 1986).

If there is at least one potential entrant, there exists a pure-strategy Nash equilibrium, at which all retailers earn zero profit by setting $w = p$, and the system operates under full efficiency.

## 2.7 Conclusions

In this chapter, we quantify the loss of efficiency from decentralizing operations in a supply chain that uses price-only contracts. In particular, we measure the Price of Anarchy, defined as the ratio of total profits between the centralized (or integrated) supply chain and the decentralized supply chain. Our bounds are IGFR distribution-free and depend on the profit margin $1 - r$.

Table 2.9 summarizes the different Prices of Anarchy bounds associated with the following supply chain configurations: push or pull inventory positioning, two or more stages, serial or assembly systems, single or multiple competing suppliers, and single or multiple competing retailers. For the sake of clarity, Table 2.9 focuses on the cases

where the downstream (resp. upstream) partner is the contract initiator when the inventory configuration is in a push (resp. pull) mode.

Table 2.9: Price of Anarchy in decentralized supply chains with price-only contracts

| Structure | $(N_C, N_S, N_R)$ | Push | Pull |
|---|---|---|---|
| 2-stage series | $(1,1,1)$ | $r^{\frac{-1}{1-r}} - r^{-1}$ | $e-1$ |
| $n$-stage series | $(1,1,1)$ | $\frac{(1-r^{\frac{-1}{n-1}})(1-r^{1-\frac{1}{1-r^{\frac{1}{n-1}}}})}{1-r}$ | $\frac{e^{n-1}-1}{n-1}$ |
| assembly | $(n,1,1)$ | $\frac{r^{1-\frac{n}{1-r}}-1}{n-1+r}$ | $e-1$ |
| procurement system[2] | $(1,n,1)$ | $1$ | $\frac{e^{\mathcal{E}^{n-1}(1)}-1}{\mathcal{E}^{n-1}(1)}$ |
| distribution system (herd) | $(1,1,n)$ | $\frac{n}{1-n(1-\frac{c}{p})}$ | $\frac{1}{n}(e^n - 1)$ |

By comparing the values of the Price of Anarchy under different supply chain network configuration, we make the following observations:

1. In simple configurations (e.g., two-stage supply chains), the Price of Anarchy is at least 1.71. Thus, the inefficiency of price-only contracts has not been overstated in the literature. This magnitude of the loss of efficiency justifies the whole stream of research on the design of more elaborate contracts and motivates future research on improving coordination in supply chains.

2. The efficiency of price-only contracts generally drops as the number of intermediaries increases (in series or assembly). However, this might not always be the case: in a pull assembly system, the Price of Anarchy is independent of the number of suppliers.

3. Introducing competition generally increases supply chain efficiency. The impact of competition on supply chain coordination is sometimes radical (push procurement), sometimes more gradual (pull procurement, push distribution). However, competition is not always beneficial and can sometimes lead to an increase in inefficiency (pull distribution).

4. In general, a pull configuration mitigates the effects of double marginalization more than a push configuration, as the inventory risk is shared among the supply

chain partners. This is especially true in an assembly system, in which the loss of efficiency might be extreme; therefore, in complex assembly systems with small profit margins, huge benefits can be reaped by moving the inventory from a push to a pull configuration. With competition however, a pull system might outperform a push system.

Table 2.9 summarizes the results obtained for the "classical" Stackelberg games, where the leader is the upstream party in the push game and the downstream party in the pull game. However, if we reverse the roles between the upstream and the downstream parties, full coordination is in general achieved, i.e., $PoA = 1$. Future research needs to address the impact of the reservation profit (assumed to be zero in this chapter) and to analyze the possible improvements through long-term relationship or renegotiation. Also, it would be interesting to quantify one party's share of profit in the worst case, as well as formalizing the impact of the coefficient of variation on supply chain efficiency.

Quantifying the efficiency of supply chains based on the Price of Anarchy generates good insights into supply chain design. Obviously, our base model is fairly idealized and can be extended in various ways. Future research needs to characterize the impact on the supply chain performance of information disclosure, retail pricing decisions, inventory sharing and allocation, and efforts to increase sales, improve forecast accuracy, or reduce costs.

# Part II

# Dynamic Traffic Assignment

# Chapter 3

# An Analytical Model for Traffic Delays and the Dynamic User Equilibrium Problem

## 3.1 Introduction

Over the last 20 years, traffic congestion has grown dramatically, resulting in high environmental costs and productivity losses. According to the Texas Transportation Institute (2003) on traffic congestion in the United States, "the average delay for every person in the 75 urban areas studied climbed from 7 hours in 1982 to 26 hours in 2001. [...] The total congestion 'invoice' for the 75 areas in 2001 came to $69.5 billion, which was the value of 3.5 billion hours of delay." All the proposed solutions (capacity increase, highway efficiency improvement, demand management) rely on an accurate prediction of traffic congestion. To this end, it has become critical: (1) to determine the travel time of a traveler and how it is affected by congestion, and (2) to understand how traffic distributes in a transportation network.

In this chapter, we derive an analytical travel-time function that integrates the traffic dynamics and the effects of shocks. Subsequently, we illustrate how this function can be employed to determine the routes that individual travelers take in a

transportation network. In particular, we assume that each user in the system minimizes the time length of his or her own trip, depending on congestion, leading to a dynamic user equilibrium (DUE) in the transportation network.

Analysis of traffic flow can be microscopic or macroscopic. Microscopic models focus on the behavior of a single vehicle, reacting to other vehicles' behavior and generally adopt a simulation approach (e.g., see Herman et al. 1959, Gazis et al. 1961, and Mahut 2000). In contrast, macroscopic models rely on the aggregate behavior of vehicles, depending on surrounding aggregate traffic conditions. Most macroscopic models are based on the theory of kinematic waves in transportation, developed by Lighthill and Whitham (1955) and Richards (1956), that models traffic flow as a compressible fluid in a pipeline. The fluid approach represents the limiting behavior of a stochastic process of a large population, and is therefore appropriate to large-scale problems, such as traffic equilibrium problems on long crowded roads. However, since macroscopic models generally assume a uniform, stationary distribution of traffic, they are generally unable to capture local interactions of vehicles, the time-dependent behavior of traffic within a queue (e.g., the "stop-and-go" phenomenon), and therefore ignore the local triggers of congestion such as traffic controls, turn permissions, opposing traffic streams, etc. Nonetheless, despite these limitations, the theory provides good estimates of the delays caused by queues and their dependence upon entrance and exit flows (Hurdle and Son 2000). In this chapter, we derive a function of congestion delays from the theory of kinematic waves.

Over the past decade, most developments based on the theory of kinematic waves have focused on modeling flow propagation rather than deriving an analytical travel-time function. Newell (1993) modeled flow propagation by using curves of the cumulative number of vehicles at the road entrance and at the road exit. Travel times correspond to the horizontal difference between these two curves. Daganzo (1994, 1995a) proposed an efficient way of propagating flow, by splitting the road into cells. (See also Velan (2000) for enhancements of this model.) His model, called the Cell Transmission Model, is not only very efficient for simulating the transportation network but can also be formulated as a linear optimization problem (Ziliaskopoulos

114

2000). While the Cell Transmission Model relies on a triangular flow-density curve, other simulation models rely on a quadratic flow-density curve (see Daganzo 1995b, and Khoo et al. 2002).

On the other hand, most traffic assignment models rely on an analytical travel-time function. In fact, traffic assignment problems (such as the DUE problem) can be formulated as nonlinear problems and be solved with standard optimization algorithms. In particular, when the travel-time function is differentiable, gradient algorithms can be used to solve traffic assignment problems (see Patriksson 1994 for a review). In addition, analytical models allow studying the convergence behavior of these traffic assignment algorithms (e.g., see Friesz et al. 1993, Ran and Boyce 1994).

However, it is unclear which analytical travel-time function should be used in these models. The most generic travel-time function depends on the number of cars, the inflow, and the outflow (see Ran and Boyce 1994). Although Daganzo (1995c) claimed that a travel-time function should depend only on the number of cars on the road, Lin and Lo (2000) pointed a paradoxical situation with such a function. On the other hand, Carey et al. (2003) introduced a generic travel-time function depending on the average flow rate, satisfying FIFO and other desirable properties. Based on the theory of kinematic waves, Perakis (2000) and Kachani and Perakis (2001) proposed polynomial and exponential travel-time functions for situations without congestion. Based on the simplified kinematic wave model proposed by Newell (1993), Kuwahara and Akamatsu (2001) derived an analytical function of the instantaneous travel time in order to solve the dynamic user equilibrium. However, the travel time they derived does not represent the actual (or experienced) travel time, unless traffic conditions remain constant. Finally, the theory of vertical queues suggests that the travel time can be decomposed into a free-flow travel time and a waiting time in a queue (e.g., see Li, Fujiwara, and Kawakami 2000). However, these travel-time functions typically do not consider other flow dynamics than the queue and disregard the spillback effects.

In this chapter, we propose a methodology for deriving a polynomial travel-time function on a single stretch of road, based on the theory of kinematic waves. This travel-time model only depends on the traffic conditions at the entrance and at the

115

exit of the road and is therefore less memory-intensive than the Cell Transmission Model. We extend the work by Perakis (2000) and Kachani and Perakis (2001) by including congestion into a travel-time function. We also complement the work by Kuwahara and Akamatsu (2001) by analyzing the dynamic user equilibrium with experienced instead of instantaneous travel times. Finally, we generalize the queuing models by integrating the flow dynamics and the spillback effects into the travel-time function. The main contributions of our model are the following:

1. We develop a methodology for determining travel times analytically that applies to both triangular and quadratic flow-density curves.

2. We introduce a travel-time function that integrates the first-order traffic dynamics, shocks, and queue spillovers.

3. From our numerical experiments, it appears that the travel-time function is consistent with the results obtained by simulation.

4. We establish that the travel-time function satisfies properties such as continuity, monotonicity, and FIFO.

5. We incorporate the travel-time model within a dynamic user equilibrium (DUE) model.

The chapter is organized as follows. In Section 3.2, we review the theory of kinematic waves. In Section 3.3, we propose a general methodology for deriving a travel-time function and illustrate it within two particular models of flow-density (quadratic and triangular). In Section 3.4, we discuss the properties of the travel-time function that we derived. In Section 3.5, we embed our travel-time function within a more general dynamic user equilibrium problem and illustrate our model through a numerical example. Finally, we conclude by outlining directions for further research.

116

## 3.2 Review of the Theory of Kinematic Waves

In this section, we review the hydrodynamic theory of traffic flow, proposed by Lighthill and Whitham (1955) and Richards (1956). From a macroscopic point of view, the flow of traffic on a stretch of a road can be modeled as the flow of a fluid in a pipeline. Accordingly, traffic flow is animated with waves, moving backwards or forwards.

Because of the dynamic nature of traffic, we work on a time-space plane. The fundamental traffic variables to describe traffic conditions on a road are:

- the flow rate, $f(x, t)$, which is the number of vehicles per hour passing location $x$ at time $t$,

- the rate of density (or concentration), $k(x, t)$, which is the number of vehicles per mile, at location $x$ at time $t$, and

- the instantaneous velocity, $v(x, t)$, which is the speed of vehicles passing location $x$ at time $t$.

Assuming a homogeneous road, these quantities are respectively bounded from above by $f^{max}$, $k^{max}$, and $v^{max}$.

Most models assume a one-to-one relation between the speed and the density, i.e., $v(x, t) = v(k(x, t))$. This relation has the additional property that when the density is zero, the speed is equal to the free-flow speed $v^{max}$, and when the density is equal to the jam density, $k = k^{max}$, the speed is zero.

Notice that the definition of the fundamental traffic variables implies that $f(x, t) = k(x, t)v(x, t)$. Therefore, the flow and the density are related through the so called *fundamental diagram*:

$$f(x, t) = k(x, t)v(k(x, t)), \quad \forall x, t. \tag{3.1}$$

Depending on the assumed speed-density relation $v(k)$, the fundamental diagram can have different shapes. In what follows, we analyze the two most common shapes, namely the quadratic and the triangular fundamental diagrams.

Greenshields (1935) modeled the vehicle velocity as a linear function of density, i.e., $v(k) = v^{max}(1 - k/k^{max})$. Based on this function, the fundamental diagram has a quadratic shape, as in Richards (1956):

$$f(x, t) = k(x, t)v^{max}(1 - k(x, t)/k^{max}). \tag{3.2}$$

Newell (1993) proposed a triangular curve with a left slope of $1/u_0$ and a right slope of $-1/w_0$. The change of slope occurs when $k = k(f^{max})$, where $k(f^{max})$ is the density associated with the road capacity.

$$f(x, t) = \begin{cases} k(x, t)/u_0 & \text{if } k(x, t) \le k(f^{max}) = f^{max}u_0, \\ (k^{max} - k(x, t))/w_0 & \text{otherwise.} \end{cases} \tag{3.3}$$

In addition to (3.1), the traffic variables are related through a conservation law, stated as the following partial differential equation:

$$\frac{\partial k(x, t)}{\partial t} + \frac{df(k)}{dk}\frac{\partial k(x, t)}{\partial x} = 0. \tag{3.4}$$

This conservation law describes the fact that on a single stretch of road, no cars are lost.

In the time-space plane, a level curve of density is the set of points $(x, t)$ such that $k(x, t)$ remains constant. Since the density uniquely determines the flow according to (3.2) or (3.3), the flow $f(k)$ and by consequent $df(k)/dk$ also remain constant along a level curve of density. Therefore, if $k(x, t)$ remains constant, (3.4) reduces to the equation of a straight line with slope $df/dk$. In other words, the level curves of density are straight lines, called characteristic lines.

Characteristic lines represent the propagation of traffic density in the time-space plane. Essentially, density propagates as a wave with speed $df/dk$. The wave speed (i.e., the slope of the characteristic line), $df/dk$, is positive when the traffic is light ($k \le k(f^{max})$) and negative when the traffic is heavy ($k \ge k(f^{max})$). Accordingly, macroscopic traffic conditions propagate forwards when traffic is light and backwards

against traffic when traffic is heavy. In particular, with a triangular fundamental diagram, the speed of a wave is $1/u_0$ in light traffic and $-1/w_0$ in heavy traffic. With a quadratic diagram, the wave speed is $v^{max}(1 - 2k(x,t)/k^{max})$, which is positive when $k(x,t) \leq k^{max}/2$, and negative otherwise.

It might seem surprising at first to have waves moving in the opposite direction of vehicles. In fact, in many dynamical systems, particles do not move in the same direction as waves. Haberman (1977) illustrates this phenomenon by a horizontal rope, attached to a tree at one side, and vertically moved by a person at the other side. Although waves propagate on the rope from the person to the tree, particles of the rope only move up and down. In transportation, a wave moves backwards when cars are decelerating: The position at which a car begins decelerating is in fact upstream of the position at which the preceding car began decelerating.

If two characteristic lines intersect, the density around the point of intersection is discontinuous. The set of points of intersection is called a shock wave and represents a sudden change in traffic conditions. Behind a bottleneck, a shock wave separates a downstream congested region from an upstream uncongested region. The shock wave propagates backwards when the queue is increasing or forwards when the queue is decreasing.

On the other hand, after an increase of capacity (e.g., when a traffic light turns green), there is a discontinuity of traffic between an upstream congested region and a downstream uncongested region. At the location of the first vehicle in the queue, traffic is discontinuous; hence, the lead vehicle is assumed to accelerate instantaneously, from zero to the free-flow speed. From the lead vehicle's location originates a fan of waves of all possible velocities. In the case of a quadratic fundamental diagram, each of the waves in the fan corresponds to a different density. Accordingly, the following vehicles in the queue accelerate gradually, as they pass through the fan of waves. In contrast, in the case of a triangular fundamental diagram, all waves in the fan have the same density $k(f^{max})$ (as their velocity is a supergradient of the fundamental diagram when the flow is at capacity); hence, the following vehicles in the queues are assumed to accelerate instantaneously, from zero to the free-flow speed.

Figure 3-1 (borrowed from Lighthill and Whitham's paper) illustrates a traffic light cycle. The straight lines are the characteristic lines, the bold line is the shock wave trajectory, and the dashed line represents a single vehicle's trajectory. When the light turns red, a queue appears behind the traffic light. The characteristic waves in the upstream uncongested region have positive slope while those in the downstream congested region (i.e., the queue) have negative slope. The two regions are separated by a shockwave, which propagates backwards as the queue grows. After the light turns green, the first car in the queue instantaneously accelerates from rest to the free-flow speed. At this point of discontinuity, characteristic lines of all intermediate slopes fan out. As the queue clears up, the shock wave moves forwards and passes through the intersection after some time. Notice that the represented vehicle crosses various waves during its trip. In particular, in the case of a quadratic fundamental diagram, all waves in the fan have different densities, and the vehicle can increase its speed only slowly, after the light turns green, as it traverses the fan of waves. More details and examples of the theory can be found in Haberman (1977).



Figure 3-1: A traffic light cycle

## 3.3 An Analytical Derivation of the Travel-Time Function

In this section, we propose a general methodology for evaluating the travel time $\tau(L, t_0)$ of a vehicle entering a homogeneous road of length $L$ at time $t_0$, based on the theory of kinematic waves. In particular, our methodology applies to the triangular and quadratic flow-density curves.

### 3.3.1 General Framework

Because of the discontinuity in density induced by shocks, the kinematic wave model may be quite hard to solve. In what follows, we introduce three assumptions that simplify the model.

**Assumption 2. Linear density.** *The second-order variation of density is locally negligible.*

Accordingly, for every period $t$, the density at location $x$ can be approximated by

$$k(x, t) = k(\xi, t) + B(\xi, t)(x - \xi), \tag{3.5}$$

where $\xi = 0$ if the traffic conditions at $(x, t)$ are light, and $\xi = L$ if they are heavy. We denote by $B(\xi, t)$ the rate of evolution of the density at time $t$ at the road entrance (if $\xi = 0$) or at the road exit (if $\xi = L$).

We focus on the rate of evolution of the density, instead of the flow, as it might seem more natural, because the value of density carries more information than the value of flow. In particular, from the fundamental diagram, a given flow can be associated with a low or high density, whereas a given density uniquely determines the flow.

Under this assumption, an analytical model of travel time of a vehicle departing at time $t_0$ and (approximately) arriving at time $t_0 + \theta \geq t_0$ (e.g., $\theta$ might be chosen to be the free flow travel time) can be based on the knowledge of only two quantities,

$B(0, t_0)$ and $B(L, t_0 + \theta)$. In contrast, an exact model would consider all different values of density $k(x, t)$ of the waves that a vehicle would cross during its trip. If the incoming flow or the exit capacity varies in a very nonlinear fashion, Assumption 2 will be violated. On the other hand, when the flow evolves smoothly, the densities of the waves that the vehicle crosses, will not be very different from each other, and the second-order variations will be negligible. In the context of traffic equilibrium, one expects that if the travel demand changes smoothly, the flow on each arc will also evolve smoothly, if the path travel-time function is continuous. This intuition has also been confirmed in our numerical example in Section 3.5.

In addition, by considering smaller road lengths, a vehicle will cross fewer waves and hence will encounter less variable traffic conditions. At the limit, one can consider road segments as small as the cells of Daganzo (1994), on which the flow remains constant. However, improved accuracy would come at the price of greater memory requirements.

In practice, parameters $B(0, t)$ and $B(L, t + \theta)$ will have to be estimated from the traffic conditions at the road entrance and exit respectively. The quantity $B(0, t + \gamma(t))$ can be estimated by $(k(0, t) - k(0, t + \gamma(t)))/L$ where $\gamma(t)$ is the time for a forward wave crossing the road entrance at time $t$ to reach the road exit, if there was no shock. In particular, $\gamma(t) = Lu_0$ in the case of a triangular flow-density curve and $\gamma(t) = L/(v^{max}(1 - 2k(L, t)/k^{max}))$ in the case of a quadratic flow-density curve. Parameter $B(L, t + \theta)$ can be estimated in a similar way.

The linear relationship of density (3.5) is not required to be time-consistent. That is, assuming a particular value for $k(x, t_0)$ at time $t_0$ does not imply that there was an actual wave with this density that crossed the road entrance some time before. Instead, Assumption 2 must be considered as a snapshot of the traffic situation from time $t_0$ onwards, for computing the travel time of a vehicle starting its trip at $t_0$.

To illustrate the implications of this assumption, let us consider the traffic light example, presented in the last section (Figure 3-1). Since we consider homogeneous roads, the traffic light must be located at the intersection of two road segments.

Let us first consider the downstream road segment. Suppose that the traffic light

122

has been red for a long time, so that there are only waves of zero density on the road. Following the theory of kinematic waves, when the traffic light turns green, the density at the entrance, $k(0,t)$, suddenly jumps to $k(f^{max})$, while the traffic density at all other points of the road remains zero, i.e., $k(x,t) = 0$ for $x \in (0, L]$. From the origin emanates a fan of waves of all possible velocities, as is illustrated in the left part of Figure 3-2.

Instead, Assumption 2 states that $k(x,t) = k(f^{max}) + B(0,t)x = k(f^{max})(L-x)/L$, since $B(0,t) = (0 - k(f^{max}))/L$. Therefore, the gradual change of density is still captured under Assumption 2 but the waves associated with these densities do not originate from the same point, as shown in the right part of Figure 3-2. Similarly, for the upstream segment, Assumption 2 captures the gradual change of density from $k^{max}$ (at the end of the queue) to $k(f^{max})$ (at the head of the queue), but the waves do not originate from the same point. Since the waves are more spread out under Assumption 2, the time to exit the queue will be underestimated and the time to accelerate on the downstream road will be overestimated.
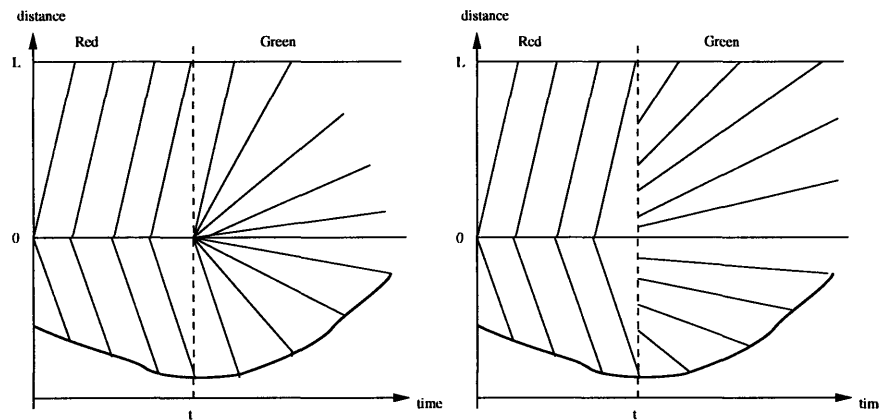


Figure 3-2: When the traffic light turns green. Comparison between the theory of kinematic waves (left) and the simplified model, under Assumption 2 (right)

.

**Assumption 3. As most one shock.** *There is at most one shock on the road, dividing an upstream uncongested region from a downstream congested region.*

123

In the original model by Lighthill and Whitham (1955) and Richards (1956), a shock may result from the focusing of two forward waves, two backward waves, or one forward and one backward wave. However, as argued by Newell (1993), only the latter type of shocks is observed in reality.

Accordingly, the road can be divided into two segments, separated by the shock wave: On the first segment, the traffic flow has a low density, whereas on the second, it has a high density. If there is no shock, the second segment has zero length, while if heavy traffic conditions back up to the road entrance, the first segment has zero length.

When the fundamental diagram is triangular, waves in a certain regime have all the same speed; hence, Newell's assumption is automatically satisfied. When the fundamental diagram is quadratic, the assumption holds if $0 \leq k(0,t) + B(0,t)L \leq k(f^{max})$ if the traffic conditions at $(x,t)$ are light, and $k(f^{max}) \leq k(L,t) - B(L,t)L \leq k^{max}$ if they are heavy (see the appendix).

**Assumption 4. Bounded variation of density.** *The variation of density is bounded by*

$$|B(\xi,t)| < (k^{max} - k(\xi,t))^2/(4Lk^{max}), \; for \; \xi = 0 \; or \; L. \qquad (3.6)$$

This assumption allows us to neglect high order terms in the travel-time function, as we will show in the proofs of Theorems 1 and 2. If the evolution of traffic flow is highly variable, we can relax (3.6) by considering smaller road lengths.

**Methodology**

From Assumption 3, the road can be decomposed into two segments, separated by the shock wave. Therefore, before a vehicle reaches the shock wave, its travel time depends only on the light traffic conditions, while after the shock wave, its travel time depends only on the heavy traffic conditions. As a result, as in Ran et al. (1997), the total travel time can be decomposed as the sum of

124

- the travel time to go from the entrance to the shock wave (under light traffic conditions), and

- the travel time to go from the shock wave to the road exit (under heavy traffic conditions).

In particular, let us consider a vehicle that starts its trip at time $t_0$, on a road of length $L$. We denote by $\tau(x, t_0)$ its travel time to reach location $x$. We assume that we know the traffic conditions (cumulative number of vehicles, density, rate of evolution of density) at the entrance at time $t_0$ and at the exit at time $t_0 + \theta$, for some $\theta \geq 0$.

**Shock Location.** A shock is a discontinuity in the traffic flow. Rather than working with flows, Newell (1993) introduced the concept of cumulative number of vehicles passing through location $x$ by time $t$, $F(x, t)$. The partial derivatives of $F(x, t)$ correspond to the density and the flow rates, i.e., $\partial F(x, t)/\partial x = -k(x, t)$ and $\partial F(x, t)/\partial t = f(x, t)$.

Along a characteristic line passing through $(x_0, t_0)$ with slope $df/dk$, the rate of evolution of $F(x, t)$ with respect to $x$ is equal to:

$$\frac{dF(x, t)}{dx} = -k(x, t) + f(x, t)(\frac{df(k)}{dk})^{-1}. \qquad (3.7)$$

Along a characteristic line, variables $k$, $f(k)$ and $df(k)/dk$ remain constant, implying that $dF/dx$ also remains constant. Therefore, the knowledge of $F(x_0, t_0)$ and $k(x_0, t_0)$ completely determines the cumulative number of vehicles at each point on the characteristic line.

If the characteristic line has positive slope, $k(x_0, t_0)$ is approximated by $k(0, t_0) + B(0, t_0)x_0$, according to (3.5). Likewise, $F(x_0, t_0) = F(0, t_0) + \int_0^{x_0} k(x, t)dx$ is approximated by $F(0, t_0) + k(0, t)x_0 + 1/2B(0, t_0)x_0^2$. As a result, the cumulative number of cars at $(x, t)$ is a function of the traffic conditions at the entrance, i.e., a function of $k(0, t_0)$, $B(0, t_0)$ and $F(0, t_0)$. In this case, we denote by $A(x, t)$ the cumulative

125

number of vehicles passing $x$ by time $t$, instead of $F(x, t)$. We use this notation in order to emphasize the reference to the cumulative number of arrivals on the road.

Similarly, if the characteristic line has negative slope, the cumulative number of vehicles at $(x, t)$ is a function of the traffic conditions at the exit of the road. In this case, we denote by $D(x, t)$ the cumulative number of vehicles passing $x$ by time $t$, instead of $F(x, t)$, in order to emphasize the reference to the cumulative number of departures from the road.

If there was no shock, all characteristic lines would never intersect, and the cumulative number of vehicles would be uniquely determined at every point. However, in the presence of shocks, two characteristic lines could potentially intersect at point $(x, t)$. In such a case, the cumulative number of vehicles at $(x, t)$ would have two different values, namely $A(x, t)$ and $D(x, t)$. Newell (1993) argued that the correct value for the cumulative number of vehicles passing $x$ by time $t$ is the minimum between $A(x, t)$ and $D(x, t)$, and that the intersection, $A(x, t) = D(x, t)$, determines the shock.

Plugging the vehicle's trajectory $(x, t_0 + \tau(x, t_0))$ into the shock wave equation, $A(x, t) = D(x, t)$, gives rise to the point $(\hat{x}, t_0 + \hat{\tau})$ at which the vehicle goes through the shock, i.e., $\hat{\tau} = \tau(\hat{x}, t_0)$.


**Travel-Time Function in Light/Heavy Traffic.** Under Assumption 3, the total travel time can be decomposed as the sum of the travel times before and after the shock.

From the fundamental diagram (3.1), the instantaneous vehicle's velocity at $(x, t)$ is the ratio between the flow and the density. Accordingly, the vehicle's trajectory $(x, t_0 + \tau(x, t_0))$ evolves as follows:

$$\frac{d\tau(x, t_0)}{dx} = \frac{1}{u(x, t_0 + \tau(x, t_0))} = \frac{k(x, t_0 + \tau(x, t_0))}{f(x, t_0 + \tau(x, t_0))}. \tag{3.8}$$

From the flow-density curve ((3.2) or (3.3)), $f(x, t_0 + \tau(x, t_0))$ can be expressed as a function of $k(x, t_0 + \tau(x, t_0))$. Therefore, the right hand side of (3.8) depends only

on density.

If the traffic conditions are light at $(x, t)$, the instantaneous velocity (3.8) can be expressed as a function of the traffic conditions at $(0, t_0)$. Along a characteristic line, the density remains constant; hence, $k(x, t_0 + \tau(x, t_0)) = k(x_0, t_0)$, where $x_0$ is the ordinate at $t_0$ of the characteristic line passing through $(x, t_0 + \tau(x, t_0))$.

Using Assumption 2, $k(x_0, t_0) = k(0, t_0) + B(0, t_0)x_0$. Plugging the density function into (3.8) gives rise to an ordinary differential equation (ODE). This ODE, together with the initial condition $\tau(0, t_0) = 0$, can be solved through a power series expansion (see Edwards and Penney 1985). Under Assumption 4, the ratio between two successive terms in the series is bounded from above by 1, and the series converges.

Similarly, if the traffic conditions are heavy at $(x, t)$, the instantaneous velocity at $(x, t)$ can be expressed as a function of the traffic conditions at $(L, t_0 + \theta)$, for some $\theta \geq 0$. Denoting by $x_\theta$ the ordinate at $t_0 + \theta$ of the characteristic line passing through $(x, t)$, we have that $k(x, t_0 + \tau(x, t_0)) = k(x_\theta, t_0 + \theta)$. Under Assumption 2, $k(x_\theta, t_0 + \theta) = k(L, t_0 + \theta) - B(L, t_0 + \theta)(L - x_\theta)$. Plugging this function of density into (3.8) gives rise to an ordinary differential equation (ODE), which can be solved with a power series solution. However, since $(x, t)$ is in the heavy traffic region, the boundary condition is now defined by the shock location, i.e., $\tau(\hat{x}, t_0) = \hat{\tau}$, where $(\hat{x}, \hat{\tau})$ solve $A(\hat{x}, t_0 + \tau(\hat{x}, t_0 + \hat{\tau})) = D(\hat{x}, t_0 + \tau(\hat{x}, t_0 + \hat{\tau}))$.

**Flow Propagation.**  A model of travel time is also a model of flow propagation. In effect, the cumulative number of vehicles at the exit at time $t + \tau(L, t)$ is equal to the cumulative number of vehicles at the entrance at time $t$, i.e., $A(0, t) = D(L, t + \tau(L, t))$.

However, because of time discretization in numerical experiments, $D(L, t)$ is not defined for the periods $t$ that do not correspond to departure times. Therefore, we have assumed that $D(L, t)$ is piecewise linear between two successive departure times, i.e., for $\lceil \tau(L, t - 1) \rceil - 1 \leq s \leq \lfloor \tau(L, t) \rfloor$,

$$D(L, t + s) = \lambda A(0, t) + (1 - \lambda)A(0, t - 1),$$

127

where $\lambda = \frac{(t+s)-(t-1+\tau(L,t-1))}{(t+\tau(L,t))-(t-1+\tau(L,t-1))}$.

At time $t$, the traffic conditions at the exit are known up to time $t-1+\tau(L,t-1)$, since they depend on the flows that entered the road at or prior to time $t-1$. Therefore, in our numerical experiments, we have chosen $\theta = \lfloor \tau(L,t-1) \rfloor - 1$.

For a one-link network, the flow at the road exit, $f(L,t+\theta)$, is set equal to the exit bottleneck capacity $Q$. For a two-link network, the exit capacity at the end of the upstream link (denoted as link 1) is time-varying since a queue may spill back from the downstream link (denoted as link 2). If the queue has uniform density, equal to $k_2(L_2,t+\theta)$, only $k_2(L_2,t+\theta)L_2$ vehicles can be present on link 2 at time $t+\theta$. Considering the difference between inflows and outflows, at most $D_2(L_2,t+\theta+1) - A_2(0,t+\theta) + k_2(L_2,t+\theta)L_2$ vehicles are allowed to enter link 2 in period $t+\theta$. As a result, the capacity at the exit of link 1 is the minimum between $Q_2$ (static capacity of link 2) and $D_2(L_2,t+\theta+1) - A_2(0,t+\theta) + k_2(L_2,t+\theta)L_2$ (dynamic capacity due to queue spillbacks). If link 1 has several downstream links or link 2 has several upstream links, the same reasoning applies, with additional care of the neighbor links. More details of the general case are presented in Section 3.5.

According to Assumption 3, the density at the entrance of the road is taken as the low density associated with the incoming flow $f(0,t)$ and the density at the exit of the road is taken as the high density associated with $f(L,t+\theta)$. As mentioned in the discussion of Assumption 2, $B(0,t+\gamma(t)) = (k(0,t) - k(0,t+\gamma(t)))/L$, where $\gamma(t) = L/(df(0,t)/dk)$. Similarly, $B(L,t+\gamma(t)) = (k(L,t+\gamma(t)) - k(L,t))/L$, where $\gamma(t) = -L/(df(L,t)/dk)$.

In what follows, we apply this general methodology to the triangular and the quadratic fundamental diagrams.

## 3.3.2  Triangular Fundamental Diagram

### Shock Location

We first show that, along the vehicle's trajectory, the cumulative number of vehicles based on the prevailing traffic conditions at the entrance, $A(x,t)$, remains constant.

Indeed, when traffic is light, the vehicle is moving at the same speed as the wave conveying the entering flow, $\tau(x, t_0) = xu_0$. Therefore,

$$
\begin{aligned}
A(x, t_0 + \tau(x, t_0)) &= A(0, t_0) + \left.\frac{dA(\xi, t_0)}{d\xi}\right|_{\xi=0} x, \\
&= A(0, t_0) + (-k(0, t_0) + f(0, t_0)u_0)x \qquad \text{from (3.7)}, \\
&= A(0, t_0) \qquad\qquad\qquad\qquad\qquad\quad \text{from (3.3)}.
\end{aligned}
$$

On the other hand, the cumulative number of vehicles depending on the traffic conditions at the exit, $D(x, t)$, varies along the vehicle trajectory. In contrast to the light traffic region, the vehicle's trajectory crosses several backward waves.

The cumulative number of vehicles in the heavy traffic region, $D(x, t)$, is expressed as a function of the the traffic conditions at $(L, t_0 + \theta)$. The ordinate at $t_0 + \theta$ of a characteristic line passing through point $(x, t_0 + \tau(x, t_0))$ is $x_\theta = x - (\theta - \tau(x, t_0))/w_0$. Because $dD/d\xi$ remains constant along a characteristic line, from (3.7), the cumulative number of vehicles based on the traffic conditions at the exit varies as follows:

$$
D(x, t_0 + \tau(x, t_0)) = D(x_\theta, t_0 + \theta) + \left.\frac{dD(\xi, t_0 + \theta)}{d\xi}\right|_{\xi=x_\theta} (x - x_\theta). \qquad (3.9)
$$

From Assumption 2, $D(x_\theta, t_0+\theta) = D(L, t_0+\theta) + k(L, t_0+\theta)(L-x_\theta) - 1/2 B(L, t_0+\theta)(L - x_\theta)^2$. Replacing $dD/d\xi$ in (3.9) with (3.7) gives rise to

$$
\begin{aligned}
D(x, t_0 + \tau(x, t_0)) &= D(L, t_0 + \theta) + k(L, t_0 + \theta)(L - x_\theta) - \frac{1}{2}B(L, t_0 + \theta)(L - x_\theta)^2 \\
&\quad + (-k(x_\theta, t_0 + \theta) - f(x_\theta, t_0 + \theta)w_0)(x - x_\theta), \\
&= D(L, t_0 + \theta) + k(L, t_0 + \theta)(L - x_\theta) - \frac{1}{2}B(L, t_0 + \theta)(L - x_\theta)^2 \\
&\quad - k^{max}(x - x_\theta),
\end{aligned}
$$

where the last equality comes from the backward wave speed definition, $w_0 f(x_\theta, t_0 + \theta) = k^{max} - k(x_\theta, t_0 + \theta)$.

At the shock location, $x_\theta = \hat{x}(1 + u_0/w_0) - \theta/w_0$, and the cumulative number of

129

vehicles becomes

$$
\begin{aligned}
D(\hat{x}, t_0 + \tau(\hat{x}, t_0)) = & \; D(L, t_0 + \theta) + k(L, t_0 + \theta)(L + \frac{\theta}{w_0}) - \frac{1}{2}B(L, t_0 + \theta)(L + \frac{\theta}{w_0})^2 \\
& -k^{max}\frac{\theta}{w_0} - k(L, t_0 + \theta)\hat{x}(1 + \frac{u_0}{w_0}) \\
& +B(L, t_0 + \theta)\hat{x}(1 + \frac{u_0}{w_0})(L + \frac{\theta}{w_0}) \\
& +k^{max}\hat{x}\frac{u_0}{w_0} - \frac{1}{2}B(L, t_0 + \theta)\hat{x}^2(1 + \frac{u_0}{w_0})^2.
\end{aligned}
$$

Equating $D(\hat{x}, t_0 + \tau(\hat{x}, t_0))$ to $A(\hat{x}, t_0 + \tau(\hat{x}, t_0)) = A(0, t_0)$ gives rise to a quadratic equation in $\hat{x}$. Considering the lower root gives the shock location when the vehicle will go through it.

In particular, when $B(L, t_0 + \theta) = 0$, the shock location is simply

$$
\hat{x} = \frac{D(L, t_0 + \theta) - A(0, t_0) + k(L, t_0 + \theta)(L + \frac{\theta}{w_0}) - k^{max}\frac{\theta}{w_0}}{k(L, t_0 + \theta)(1 + \frac{u_0}{w_0}) - k^{max}\frac{u_0}{w_0}}. \tag{3.10}
$$

**Special Case:** $\theta = 0$. We consider the special case where the available traffic information at the exit is at time $t_0$. Our goal is to justify why only $D(L, t_0) - A(0, t_0 - 1) + k(L, t_0)L$ vehicles are authorized to enter the road, when a queue spills back.

By definition, $\hat{x}$ is the shock location on the trajectory of a vehicle that enters the road at time $t_0$. For this reason, when a queue spills back onto upstream links, the shock location $\hat{x}$ must be equal to zero, since the entire road is covered by the queue. Therefore, the incoming flow will be constrained to maintain $\hat{x} \geq 0$. If $\theta = 0$, the shock location (3.10) is identified as $\hat{x} = (D(L, t_0) - A(0, t_0) + k(L, t_0)L)/(k(L, t_0)(1 + \frac{u_0}{w_0}) - k^{max}\frac{u_0}{w_0})$. The denominator is always positive since, from (3.3), $k(L, t_0) = k^{max} - f(L, t_0)w_0$ (high density value), and $(k^{max} - f(L, t_0)w_0) - f(L, t_0)u_0 \geq 0$. Therefore, $\hat{x} \geq 0$ if the numerator is positive, that is, $D(L, t_0) - A(0, t_0) + k(L, t_0)L \geq 0$.

130

**Travel-Time Function**

Because vehicle's speed is constant in light traffic, the travel time to go to the shock location is equal to $u_0\hat{x}$.

On the other hand, the travel time to go from $\hat{x}$ to the road exit is obtained by plugging $f(x, t_0 + \tau(x, t_0)) = 1/w_0(k^{max} - k(x, t_0 + \tau(x, t_0)))$ into (3.8).

Since the characteristic line passing through $(x, t_0 + \tau(x, t_0))$ intersects the time line $t_0 + \theta$ at $x_\theta = x - (\theta - \tau(x, t_0))/w_0$, and the density remains constant along a characteristic line, $k(x, t_0 + \tau(x, t_0)) = k(x_\theta, t_0 + \theta)$. Moreover, from (3.5), $k(x_\theta, t_0 + \theta) = k(L, t_0 + \theta) - B(L, t_0 + \theta)(L - x_\theta)$. Updating (3.8) gives rise to

$$\frac{d\tau(x, t_0)}{dx} = \frac{(k(L, t_0 + \theta) - B(L, t_0 + \theta)(L - x_\theta))w_0}{k^{max} - (k(L, t_0 + \theta) - B(L, t_0 + \theta)(L - x_\theta))}. \tag{3.11}$$

As shown in Theorem 19, solving this differential equation gives rise to a closed form solution of the travel-time function. The static parameters of the obtained travel time are the positive and negative wave paces, $u_0$ and $w_0$, and the road length $L$. On the other hand, the travel time also depends on the dynamic evolution of traffic, namely the cumulative number of vehicles at the road entrance and at the road exit, $A(0, t_0)$ and $D(L, t_0 + \theta)$, the associated densities, $k(0, t_0)$ and $k(L, t_0 + \theta)$, and the rate of evolution of the latter with respect to distance, $B(L, t_0 + \theta)$. Notice that the travel-time function is independent of the rate of evolution of the density at the entrance, $B(0, t_0)$.

In order to prove Theorem 19, we need a technical assumption about the value of $\theta$. Specifically, we assume that $\theta$ is greater than the free flow travel time up to the shock location and less than the road travel time with a shock.

**Theorem 19.** *Under Assumptions 2 and 4, if $u_0\hat{x} + (L - \hat{x})w_0 k(L, t_0 + \theta)/(k^{max} - k(L, t_0 + \theta)) \geq \theta \geq \hat{x}u_0$, the triangular diagram (3.3) gives rise to the following travel time for a vehicle entering the road at time $t_0$:*

$$\tau(L, t_0) = \hat{x}u_0 + \frac{k(L, t_0 + \theta)w_0 - B(L, t_0 + \theta)(w_0(L - \hat{x}) + \theta - u_0\hat{x})}{k^{max} - k(L, t_0 + \theta) + B(L, t_0 + \theta)(L - \hat{x} + (\theta - u_0\hat{x})/w_0)}(L - \hat{x})$$
$$+ O((L - \hat{x})w_0), \tag{3.12}$$

131

*where $\hat{x}$ is defined by (3.10).*

## Numerical Comparison

Daganzo (1994) proposed a discrete simulation method of flow propagation, called the Cell Transmission Model, in the case of a triangular flow-density relation. His algorithm is available through a user-friendly software program, called Netcell (Cayford et al. 1997). Netcell builds the curves of cumulative number of vehicles at the entrance and at the exit of the road. The travel time of a vehicle corresponds to the horizontal difference between the two curves.

Figure 3-3 displays the analytical travel time (3.12) and the travel time obtained with Netcell, as a function of the departure time. The entering flow evolves quadratically with time, peaking at 1600 vehicles/hour after 1/2 hour: $f(0, t) = 1600 - 6400(t/3600 - 0.5)^2$ vehicles/hour, for $t = 1, ..., 3600$ seconds. The flow-density relation is asymmetric triangular, with $u_0 = 1/40$ hour/mile, $w_0 = 1/10$ hour/mile, and $k^{max} = 200$ vehicles/mile; hence the road capacity is $f^{max} = 1600$ vehicles/hour. The road has a length of 4 miles and has a bottleneck at its end authorizing only 1400 vehicles/hour to exit the road.

The travel times determine the flow propagation, and the flow propagation determine future travel times (by building the curve $D(L, t)$). In this example, the exit bottleneck capacity is fixed, and $B(L, t)$ is zero for all $t$.

As illustrated in Figure 3-3, for each of the 3600 departure times, the analytical travel-time function practically coincides with the Cell Transmission Model. We performed many different numerical tests (with or without shock, low/high entering flow), and all displayed the same behavior as the example we study in Figure 3-3.

The next example investigates a queue spillback situation in a two-link network. The lengths of the links are 3 miles for the first one (upstream) and 1 mile for the second (downstream). The flow incoming to link 1 is the same as before, but there is now an incoming flow at the entrance of the second link, piecewise linear: $(300 - 600|t/3600 - 0.5|)$ vehicles/hour for $t = 1, ..., 3600$. The two links have the same characteristics as above ($u_0 = 1/40$ hour/mile, $w_0 = 1/10$ hour/mile and $k^{max} =$

Figure 3-3: Travel time with an exit bottleneck. Comparison between the analytical model (3.12) and the Cell Transmission Model.

200 vehicles/mile), and there is a bottleneck at the end of the second link of 1400 vehicles/hour.

Figure 3-4 displays the travel time of the downstream link (bottom), the travel time of the upstream link (middle), and the path travel time (top), as a function of the journey departure time, computed with the analytical model (3.12) and the Cell Transmission Model. As shown on Figure 3-4, a queue appears behind the bottleneck, at the end of link 2, propagates backwards on link 2 and then spills back into link 1. At this point, the travel time on link 1 starts increasing, while the travel time on link 2 remains constant, since the size and the density of the queue remain constant. Notice that, because the flow incoming to link 2 is time-varying, the capacity at the exit of link 1 is also time-varying, and $B(L, t + \theta)$ is nonzero. As shown in the figure, both the analytical model and the simulation model give rise to similar travel times.

The final example examines a queue dissipation situation. We consider the same two-link network topology as before, with the same quadratic incoming flow rate. At the end of the first link, after 3 miles, an incident occurs between periods 1100 and 1500, reducing the exit capacity from 1600 vehicles/hour to 100 vehicles/hour. Figure 3-5 displays the evolution of the travel times on link 1 (middle), link 2 (bottom), and the total path (top), using the analytical function and the Cell Transmission Model.

133

Figure 3-4: Travel time with queue spillback. Comparison between the analytical model (3.12) and the Cell Transmission Model.

Notice that the travel times on link 2 obtained with both methods are equal and correspond to the free-flow travel time, despite the approximation of wave velocities shown in Figure 3-2. In fact, with a triangular fundamental diagram, all waves that have a pace between $u_0$ and $-w_0$ have density $k(f^{max})$, since their slopes correspond to the supergradients of $f^{max}$ with respect to $k$; hence, the speed of a vehicle that crosses the fan of waves is $1/u_0$, with or without the linearity assumption. Therefore, Assumption 2 does not affect the travel time in this case.

In summary, in this subsection, we analyzed the traffic delays experienced by a traveler on a single stretch of road, when the fundamental diagram is triangular. In particular, under Assumptions 2 and 4,

- we proposed a closed-form solution of the shock location,

- we derived a closed-form solution of the travel time, capturing the traffic dynamics and spillback effects, based on the solution of a single ordinary differential equation, and

- we compared the analytical travel time with Daganzo's Cell Transmission Model, revealing that both models lead to similar results.

134

Figure 3-5: Travel time with temporary incidents. Comparison between the analytical model (3.12) and the Cell Transmission Model.

### 3.3.3 Quadratic Fundamental Diagram

In this section, we derive a travel-time function when the fundamental diagram is quadratic. As we will see, the same methodology applies and the travel-time function that we obtain is also very close to the numerical simulations.

**Travel Time without Shock**

From the quadratic fundamental diagram, we substitute $f(x, t_0 + \tau(x, t_0)) = k(x, t_0 + \tau(x, t_0))v^{max}(1 - k(t_0 + \tau(x, t_0))/k^{max})$ into (3.8).

A characteristic line passing through $(x, t_0 + \tau(x, t_0))$ intersects the time origin $t_0$ at

$$x_0 = \frac{x - v^{max}\tau(x, t_0)(1 - 2k(0, t_0)/k^{max})}{1 - 2v^{max}\tau(x, t_0)B(0, t_0)/k^{max}}.$$

Because the density remains constant along a characteristic line, $k(x, t_0 + \tau(x, t_0)) = k(x_0, t_0)$. Under Assumption 2, $k(x_0, t_0) = k(0, t_0) + B(0, t_0)x_0$. Substituting this linear function of density into (3.8) gives rise to an ODE. This ODE, together with the initial condition $\tau(0, t_0) = 0$, can be solved with a power series solution, giving rise

135

to the following analytical function of travel time in light traffic:

$$\tau(x, t_0) = \frac{x}{v^{max}(1 - k(0, t_0)/k^{max})} + \frac{1}{2}\frac{B(0, t_0)k(0, t_0)x^2}{(k^{max})^2 v^{max}(1 - k(0, t_0)/k^{max})^3} + O(\frac{L}{v^{max}}) \quad (3.13)$$

This power series converges under Assumption 4 since the ratio of two successive terms is bounded above by $|2B(0, t)L(k^{max}/(k^{max} - k(0, t))^2)|$, which is less than 1 under Assumption 4. Notice that if there is no variation in density, i.e., $B(0, t) = 0$, the travel time (3.13) reduces to the ratio of distance over speed.

**Shock Location**

Unlike with the triangular fundamental diagram, the vehicle's trajectory crosses several waves in light traffic when the fundamental diagram is quadratic. Nonetheless, the cumulative number of vehicles depending on the prevailing conditions at the road entrance, $A(x, t)$, remains constant along the vehicle's trajectory. Let us denote by $x_0$ the ordinate at $t_0$ of the characteristic line passing through $(x, t_0 + \tau(x, t_0))$. From (3.7), the cumulative number of vehicles along the vehicle's trajectory is equal to:

$$
\begin{aligned}
A(x, t_0 + \tau(x, t_0)) &= A(x_0, t_0) + \frac{dA(\xi, t_0)}{d\xi}\bigg|_{\xi=x_0} (x - x_0) \\
&= A(x_0, t_0) + (-k(x_0, t_0) + \frac{k(x_0, t_0)(1 - k(x_0, t_0)/k^{max})}{1 - 2k(x_0, t_0)/k^{max}})(x - x_0) \\
&= A(0, t_0) - \int_0^{x_0} k(x, t_0)dx + \\
&\quad (-k(x_0, t_0) + \frac{k(x_0, t_0)(1 - k(x_0, t_0)/k^{max})}{1 - 2k(x_0, t_0)/k^{max}})(x - x_0).
\end{aligned}
$$

Replacing $x_0$ by (3.13), and $\tau(x, t_0)$ by the $n$-th order expansion of (3.13) into the above expression of $A(x, t_0 + \tau(x, t_0))$, and finally taking the $n$-th order Taylor series expansion around $x = 0$, we obtain that

$$A(x, t_0 + \tau(x, t_0)) = A(0, t_0) + O(x^n).$$

Therefore, similarly to triangular case, the cumulative number of vehicle remains

136

constant along the vehicle's trajectory, despite the fact that the vehicle crosses waves of different densities.

On the other hand, $D(x, t)$ is expressed as a function of the traffic conditions at $(L, t_0 + \theta)$. Let $x_\theta$ be the ordinate at $t_0 + \theta$ of the characteristic line passing through $(x, t_0 + \tau(x, t_0))$. To simplify the computations, we assume that the density is constant in the heavy traffic region. Therefore, the characteristic line has a slope of $v^{max}(1 - 2k(L, t_0 + \theta)/k^{max})$ and $x_\theta = x + (\theta - \tau(x, t_0))v^{max}(1 - 2k(L, t_0 + \theta)/k^{max})$. From (3.7),

$$
\begin{aligned}
D(x, t_0 + \tau(x, t_0)) &= D(x_\theta, t_0 + \theta) + \frac{dD(\xi, t_0 + \theta)}{d\xi}\bigg|_{\xi=x_\theta} (x - x_\theta) \\
&= D(L, t_0 + \theta) + k(L, t_0 + \theta)(L - x_\theta) \\
&\quad + (-k(L, t_0 + \theta) + \frac{f(L, t_0 + \theta)}{v^{max}(1 - 2k(L, t_0 + \theta)/k^{max})})(x - x_\theta).
\end{aligned}
$$

At the shock location, the vehicle's travel time $\tau(\hat{x}, t_0)$, denoted by $\hat{\tau}$, is defined as in (3.13), with $\hat{x}$ instead of $x$. Plugging $\tau(\hat{x}, t_0)$ into $x_\theta$ defines the cumulative number of vehicles at the shock location as follows:

$$
D(\hat{x}, t_0 + \tau(\hat{x}, t_0)) = D(L, t_0 + \theta) + k(L, t_0 + \theta)(L - \hat{x})\frac{k(L, t_0 + \theta) - k(0, t_0)}{k^{max} - k(0, t_0)}.
$$

A shock occurs when $A(\hat{x}, t_0 + \tau(\hat{x}, t_0)) = D(\hat{x}, t_0 + \tau(\hat{x}, t_0))$, i.e., at location

$$
\hat{x} = L - (A(0, t_0) - D(L, t_0 + \theta))\frac{k^{max} - k(0, t_0)}{k(L, t_0 + \theta)(k(L, t_0 + \theta) - k(0, t_0))}. \tag{3.14}
$$

**Travel-Time Function with Shocks**

In this subsection, we integrate the effects of a shock into the travel-time function. As with the triangular flow-density curve, the resulting travel-time function has a closed form and depends on static road features (maximum speed $v^{max}$, jam density $k^{max}$, and length $L$) and on the dynamic evolution of the cumulative number of vehicles at the entrance and at the exit, $A(0, t_0)$ and $D(L, t_0 + \theta)$, the associated densities, $k(0, t_0)$ and $k(L, t_0 + \theta)$, and their rate of evolution, $B(0, t_0)$ and $B(L, t_0 + \theta)$. Similarly to

137

the triangular case, we need a technical assumption on the value of $\theta$ to prove the theorem. Specifically, we assume that it is greater than the free flow travel time up to the shock and less than the road travel time with the shock.

**Theorem 20.** *Under Assumptions 2-4, if $\hat{\tau} + (L - \hat{x})/(v^{max}(1 - k(L, t_0 + \theta)/k^{max})) \geq \theta \geq \hat{\tau}$, the quadratic diagram (3.2) gives rise to the following travel time for a vehicle entering the road at time $t_0$:*

$$\tau(L, t_0) = \hat{\tau} + \frac{(k^{max} + 2v^{max}(\theta - \hat{\tau})B(L, t_0 + \theta))(L - \hat{x})}{v^{max}(k^{max} - k(L, t_0 + \theta) + B(L, t_0 + \theta)(L - \hat{x} + v^{max}(\theta - \hat{\tau})))}$$
$$+ O(\frac{L - \hat{x}}{v^{max}}), \tag{3.15}$$

*where $\hat{x}$ is defined in (3.14), $\hat{\tau} = \tau(\hat{x}, t_0)$ according to (3.13).*

## Numerical Comparison

We compared the analytical travel time (3.15) with several travel-time functions that were proposed in the literature, all based on the theory of kinematic waves, obtained either via simulation or via an analytical derivation.

In the case of a quadratic fundamental diagram, Kachani and Perakis (2001) derived two functions (polynomial and exponential) of travel time. On the other hand, Khoo et al. (2002) simulated the flow propagation along Godunov's scheme for calculating the numerical flux. They used two different methods, called hyperbolic PDE and iterative position update methods. Finally, Daganzo (1995b) approximated the kinematic wave model with finite difference equations (FDE) and simulated flow propagation accordingly. The travel time is the horizontal difference between the cumulative departures and the cumulative arrivals on the road.

In Figure 3-6, we consider a quadratic incoming flow rate $f(0, t) = 1600 - 6400(t/3600 - 0.5)^2$ vehicles/hour, for $t = 1, ..., 3600$ seconds. The road has a length of 4 miles, a maximum speed $v^{max} = 40$ miles/hour, a maximum density $k^{max} = 200$ vehicles/mile, and hence, a capacity of $f^{max} = 2000$ vehicles/hour. There is no exit bottleneck. The x-marks represent the analytical travel time derived in Theorem

2. The three top lines are the travel times obtained by simulation by Khoo et al. (KLPP) and Daganzo (D), and the two bottom lines are the analytical functions proposed by Kachani and Perakis (KP). As illustrated in the figure, the analytical travel-time function (3.15) is very close to what was obtained from the simulations.



Figure 3-6: Travel time under light traffic, with a quadratic incoming flow rate. Comparison among the analytical model (3.15), the analytical models by Kachani and Perakis (2001), the simulations by Khoo et al. (2003) and the Finite Difference Equation simulations by Daganzo (1995b).

The second example investigates the evolution of travel times, in the presence of an exit bottleneck of 1500 vehicles/hour. In this example, we only compared the analytical travel time (3.15) with Daganzo's FDE model. In fact, the simulation results obtained by Khoo et al. were not available to us in this case, and the analytical travel-time functions proposed by Kachani and Perakis are only valid under light traffic. As illustrated in Figure 3-7, both the analytical model (3.15) and the FDE simulation model lead to similar travel times. We obtained comparable results with larger incoming flows.

Because the bottleneck capacity remains constant over time, $B(L, t + \theta) = 0$. On the other hand, the incoming flow is time-varying, and $B(0, t)$ is different than zero. In contrast, most models of vertical queues ignore the flow dynamics, setting $B(0, t)$ to zero. As illustrated in Figure 3-7, although a model of vertical queue gives a good

Figure 3-7: Travel time with an exit bottleneck. Comparison between the analytical model (3.15) and the Finite Difference Equation simulations by Daganzo (1995b).

approximation of the simulated travel time, it is not as accurate as the analytical model (3.15).

The last example analyzes a situation with a queue dissipation. We consider the same quadratic entering flow and the same road characteristics as before. We assume that an incident occurs at mile 3 after 1100 seconds, causing a temporary bottleneck of 100 vehicles/hour (in lieu of 2000 vehicles/hour) until time 1500. Figure 3-8 displays the travel times on the first 3 miles (middle), on the last mile (bottom), and on the total link (top), as a function of the departure time, obtained with the analytical model (3.15) and the FDE simulation model. The parameter $B(0, t)$ was constrained to be no larger than (3.6), in order to ensure convergence of the series expansion of (3.15). The travel time on the first 3 miles computed with (3.15) is slightly larger than the travel time obtained with simulation. In fact, the travel time (3.15) is evaluated in reference to the exit capacity at time $t + \theta = t - 1 + \lfloor \tau(t - 1) \rfloor$ and identifies the capacity increase at the end of the incident with a one-period lag.

On the other hand, the travel time on the last mile obtained with (3.15) slightly overestimates the one obtained by simulation, just after the removal of the bottleneck. In fact, Assumption 2 approximates the fan of waves appearing just after the bottleneck removal (see Figure 3-2). Since low-density waves are located farther from the
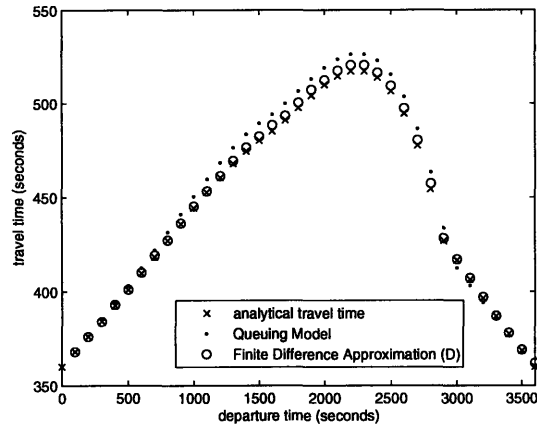
140

Figure 3-8: Travel time with an incident. Comparison between the analytical model (3.15) and the Finite Difference Equation simulations by Daganzo (1995b).

entrance under Assumption 2, the analytical travel-time function (3.15) overestimates the simulated travel time on the downstream road segment. On the other hand, after the queue has cleared up, the analytical travel time (3.15) underestimates the one obtained by simulation. In fact, the analytical travel time, based on Assumption 3, ignores the shock occurring between low-density upstream flow and and higher, but still low, density downstream flow.

In summary, in this subsection, we analyzed the traffic delays experienced by a traveler on a single stretch of road, in the case of a quadratic fundamental diagram. In particular, under Assumptions 2, 3 and 4,

- we proposed a closed-form solution of the shock location, assuming constant density in the heavy traffic region,

- we derived a closed-form solution of the travel time, capturing the traffic dynamics, based on the solution of a single ordinary differential equation,

- we highlighted the quality of prediction of the analytical travel time in comparison to simulations.

141

# 3.4 Properties of the Travel-Time Function

In this section, we analyze some common properties of the two travel-time functions that we derived in Section 3.3. First, we consider two particular cases, namely when $\theta$ corresponds to the free flow travel time and when $t + \theta$ corresponds to the departure time of the flow that arrived one period earlier. Notice that these two particular cases satisfy the conditions of Theorems 1 and 2. Second, we prove that the travel-time function is continuous, monotone and satisfies the FIFO property under certain conditions.

## 3.4.1 Two Particular Cases

### $\theta$ = Free Flow Travel Time

When $\theta$ corresponds to the free flow travel time and the dynamic effects are neglected, the travel-time functions (3.12) and (3.15) are very intuitive, as shown in the next Corollary.

**Corollary 3.** *When $\theta$ equals the free flow travel time (* $Lu_0$ *under a triangular fundamental diagram, and* $L/(v^{max}(1 - k(0,t)/k^{max})$ *under a quadratic fundamental diagram), and when* $B(0,t) = B(L, t + \theta) = 0$*, the travel-time functions (3.12) and (3.15) simplify to*

$$\tau(L,t) = \theta + \frac{A(0,t) - D(L, t + \theta)}{f(L, t + \theta)}. \tag{3.16}$$

Under the assumptions of Corollary 1, the travel-time function can be decomposed as the sum of:

- the time to travel a distance $L$ at the free flow speed, $\theta$, and

- the time to wait in a queue for available downstream capacity, $\frac{A(0,t) - D(L, t+\theta)}{f(L, t+\theta)}$.

Such a travel-time model is very similar to the traditional point queue models (e.g., see Li, Fujiwara and Kawakami 2000). However, queuing models usually assume that

142

the exit capacity is fixed, i.e., $f(L, t+\theta) = Q$. In this particular case, despite the fact that $B(0, t) = B(L, t + \theta) = 0$, the exit capacity might be different for each period of incoming flow.

### $t + \theta =$ Departure Time of the Flow that Arrived in Period $t - 1$

When $t + \theta$ corresponds to the departure time of the flow that arrived in period $t - 1$, and when $B(0, t) = B(L, t + \theta) = 0$, the expression for the travel-time function greatly simplifies and is insightful. Notice that, when we compute the travel time of the incoming flow at time $t$, time $t + \theta = t - 1 + \tau(t - 1)$ corresponds to the last period for which we have traffic information at the exit.

**Corollary 4.** *When* $\theta = \tau(t - 1) - 1$, *and* $B(0, t) = B(L, t + \theta) = 0$, *the travel-time functions (3.12) and (3.15) simplify to*

$$\tau(L, t) = \tau(L, t - 1) - 1 + \frac{f(0, t)}{f(L, t + \theta)}. \tag{3.17}$$

The corollary follows from the fact that, when $\tau(L, t) = \tau(L, t-1) - 1$, $D(L, t+\theta) = A(0, t - 1)$, and $f(0, t) = A(0, t) - A(0, t - 1)$.

Therefore, the departure time of a flow incoming at time $t$, i.e., $t + \tau(L, t)$, is the sum of:

- the departure time of the flow that arrived one period before, and

- the time to wait for the incoming flow to exit, at the rate defined by the exit capacity at time $t + \theta$.

Although both particular cases assume $B(0, t) = B(L, t + \theta) = 0$, the second case is actually more accurate than the first one. Figure 3-9 illustrates a situation where a vehicle departing at time $t$ is in a queue, assuming a triangular fundamental diagram. The horizontal distance between the curves of cumulative number of vehicles correspond to the travel time. The slope of the cumulative number of vehicles at the exit is changing over time, because the exit capacity is time-varying (e.g., due to the changing traffic on neighbor roads).

143

The first case considers $\theta = Lu_0$. The travel time is the sum of $\theta$ and the size of the queue, $A(0,t) - D(L, t + \theta)$, divided by the exit capacity $f(L, t + Lu_0)$. Since the difference between the experienced travel time and the free-flow travel time can be large, approximating the average capacity between $t + \theta$ and $t + \tau(L, t)$ with $f(L, t + \theta)$ might lead to a poor estimation of the travel time.

The second case considers $\theta = \tau(L, t-1) - 1$. The travel time is the sum of $\theta$ and the incoming flow $f(0, t)$ divided by the exit capacity $f(L, t - 1 + \tau(L, t - 1))$. Since the difference between the departure times of two successive vehicles, $\tau(L, t) + 1 - \tau(L, t-1)$, typically covers few periods, the average capacity between $t - 1 + \tau(L, t-1)$ and $t + \tau(L, t)$ will be fairly accurately approximated with $f(L, t - 1 + \tau(L, t - 1))$ (in the figure, they are equal). As a result, the approximation of the travel time with $\theta = \tau(t - 1) - 1$ will in general be more accurate than that with $\theta$ equal to the free-flow travel time.



Figure 3-9: Evolution of the cumulative number of cars at the entrance and at the exit.

In fact, even if $B(L, t + \theta) = 0$, the travel-time function (3.17) integrates the changes of exit capacity that occurred prior to $t - 1 + \tau(L, t-1)$, because it is defined recursively. The rate of change of the exit capacity $B(L, t + \theta)$ only matters for describing changes in traffic conditions between $t - 1 + \tau(L, t - 1)$ and $t + \tau(L, t)$, which typically covers few periods.

In the example of Figure 3-4, $B(L, t + \theta)$ was assumed to be nonzero to capture the dynamics of the entering traffic onto link 2. However, considering a nonzero value of $B(L, t+\theta)$ led to some numerical instability if the travel time (3.12) was truncated to low orders. In particular, truncating the travel-time function to low orders gives rise to a chain of errors: inaccurate travel times lead to inaccurate flow propagation, and vice versa. To ensure numerical stability, we had to consider a 15-th order series expansion of the travel time (3.12). In contrast, taking $B(L, t + \theta)$ equal to zero is very stable numerically. Moreover, the travel times computed with (3.17) are still very accurate: in the example of Figure 3-4, the maximum relative difference between the 15-th order series expansion of (3.12) and (3.17) was 0.25%.

This observation also holds in a less restrictive environment, when $B(0, t)$ is not required to be zero (under a quadratic diagram).

## 3.4.2 Continuity and Monotonicity

Continuity and strict monotonicity of the travel-time function are desirable properties for solving the dynamic traffic equilibrium problem. In particular, the DUE problem can be formulated as a variational inequality over a compact set (see Nagurney 1993). If the travel-time function is continuous, the variational inequality has a solution. If the travel-time function is strictly monotone, the variational inequality has at most one solution.

In what follows, we assume that $\theta$ is smaller than or equal to $\tau(L, t - 1) - 1$ and is independent of the incoming flow. Accordingly, the traffic quantities at the road exit (i.e., $D(L, t + \theta), f(L, t + \theta)$) are independent of the entering flow. For instance, the conditions are met if $\theta = \tau(L, t - 1) - 1$ or if $\theta$ corresponds to the free-flow travel time under a triangular fundamental diagram. However, with a quadratic fundamental diagram, when $\theta$ corresponds to the free-flow travel time, it depends on the incoming flow, as $\theta = L/(v^{max}(1 - k(0, t)/k^{max}))$.

**Theorem 21.** *If $\theta \leq \tau(L, t - 1) - 1$ and is independent of the incoming flow, the travel-time functions (3.12) and (3.15) are continuous.*

145

**Theorem 22.** *If $\theta \leq \tau(L, t - 1) - 1$ and is independent of the incoming flow, there exists a range of values for $B(0, t)$ and $B(L, t + \theta)$ around zero, for which the travel-time functions (3.12) and (3.15) are monotone and strictly monotone functions of flow respectively.*

Incidentally, the Jacobian matrix is lower triangular, consistently with the principle of causality (Daganzo 1995b) that states that the travel time can only be affected by current or previous flows.

### 3.4.3    First-In-First-Out (FIFO)

The First-In-First-Out (FIFO) condition guarantees that a vehicle that departed later cannot arrive earlier. The FIFO property gives consistency to the model, especially when one computes the path travel times as the sum of the link travel times. However, despite its nice analytical and computational properties, FIFO is not always verified in practice (e.g., in multilane intersections).

In the next theorem, we prove that the travel-time functions (3.12) and (3.15) satisfy the FIFO property, for the particular case where $\theta = \tau(t - 1) - 1$ and $B(L, t + \theta) = 0$.

**Theorem 23.** *If $\theta = \tau(t - 1) - 1$ and $B(L, t + \theta) = 0$, there exists a range of values for $B(0, t)$ around 0 such that the travel-time functions (3.12) and (3.15) satisfy the FIFO property.*

The FIFO property can be shown in more general cases, but we omit the proof for the sake of brevity.

## 3.5    Integration within a DUE Problem

### 3.5.1    Path Formulation of the DUE

The model for traffic delays can be embedded within a DUE setting. In this subsection, we propose a standard path formulation of a discrete-time DUE, defined on a time-expanded transportation network (see Ran and Boyce 1994).

146

Let us consider a network with a set of arcs $A$, a set of paths $P$, over a time horizon of $T$ periods. Let $W$ be the set of all origin-destination (OD) pairs. For a particular OD pair $w \in W$, let $P^w \subset P$ be the set of paths linking this origin to this destination, and let $d^w(t)$ be the given demand for period $t$. A DUE based on Wardrop's first principle (Wardrop 1952) satisfies the following property: If the flow on path $p \in P^w$ at time $t$, $f_p(t)$, is positive, the associated travel time $\tau_p(t)$ is minimal. Let $\pi^w(t)$ be the smallest travel time for the OD pair $w$ at time $t$. Mathematically, Wardrop's first principle can be formulated as follows:

$$\tau_p(t) \begin{cases} = \pi^w(t) & \text{if } f_p(t) > 0 \\ \geq \pi^w(t) & \text{if } f_p(t) = 0 \end{cases} \quad \forall p \in P^w, \forall w \in W. \qquad (3.18)$$

Stated differently, according to Wardrop's first principle, each traveler non-cooperatively seeks to minimize his/her travel time. If path $p$ consists of the $m$ arcs $\{a_1, a_2, ..., a_m\}$, the path travel time for a vehicle starting its trip at time $t$, $\tau_p(t)$, is defined recursively as

$$\tau_p(t) = \tau_{a_1}(t) + \tau_{a_2}(t + \tau_{a_1}(t)) + ... + \tau_{a_m}(t + \tau_{a_1}(t) + \tau_{a_2}(t + \tau_{a_1}(t)) + ...), \quad (3.19)$$

where $\tau_{a_i}$, $i = 1, ..., m$, are the arc travel-time functions defined in (3.12) or (3.15), depending upon the assumed curve in the fundamental diagram.

In the following, we denote by $\mathbf{f}$ the vector of path flows for all periods, i.e.,

$$\mathbf{f} = (f_1(1), f_2(1), ..., f_{|P|}(1), f_1(2), ..., f_1(T), ..., f_{|P|}(T)),$$

associated with a vector of path travel times $\tau(\mathbf{f})$. A vector of path flows is feasible if it is nonnegative and if it satisfies the demand, that is, if it belongs to the following polyhedron:

$$\mathcal{K} = \{\mathbf{f} : \sum_{p \in P^w} f_p(t) = d^w(t), \forall w \in W, \mathbf{f} \geq \mathbf{0}\}.$$

Wardrop's first principle can be expressed as the following variational inequality

147

(see Nagurney 1993):

$$\tau^*(\mathbf{f})'[\mathbf{f} - \mathbf{f}^*] \geq 0, \quad \forall \mathbf{f} \in \mathcal{K}. \tag{3.20}$$

However, there is no guarantee for the existence and the uniqueness of a solution to this variational inequality in a discrete-time setting. Even if the arc travel-time function (3.12) is continuous (see Theorem 21), the path travel times can be discontinuous, since (3.19) involves discrete time indices. This problem would not occur with a continuous-time formulation; however, the variational inequality would have infinite dimension.

On the other hand, even if the Jacobian matrix of the arc travel times (3.15) is positive definite (see Theorem 22), the Jacobian matrix of the path travel times is much more complex to analyze, since it relies on the path flow pattern (if a path flow increases, it may affect the travel time of another path flow that shares some arcs in common) and the network structure (the capacity at the exit of an arc, $Q(t)$, depends on the flow exiting adjacent arcs). In particular, the Jacobian matrix of the path travel times is not necessarily lower triangular, since the travel time of a flow starting its trip at time $t$ might be affected by a flow on another path departing some time later. Further research is needed to characterize the structure of the Jacobian matrix of the path travel times. Alternatively, one could consider a formulation of the variational inequality in terms of arcs.

### 3.5.2 Path Flow Disaggregation

Formulating (3.20) relies on mapping path flows into path travel times. However, in Section 3.3, we only mapped arc flows into link travel times. In order to use our previous results, we need to develop a procedure for disaggregating path flows into link flows. Once we know the link flows, we can easily compute the link travel times with (3.12) or (3.15) and obtain the path travel times through (3.19).

We first outline the general procedure for the path flow propagation. Working with one path at a time, we start propagating the path flow from the path origin and

continue propagating it along the path, moving forward in time. Specifically, if the path consists of arcs $a_1, a_2, ...$, we first compute the travel time on $a_1$ of the path flow starting its trip at time $t$, $\tau_{a_1}(L_{a_1}, t)$. Then, we propagate the path flow along $a_1$, and we compute its travel time on $a_2$ given that it is at the entrance of $a_2$ at time $t + \tau_{a_1}(L_{a_1}, t)$, and so on.

However, in order to compute the link travel time, one needs the total incoming flow, and not only the flow related to that particular path. Therefore, we stop propagating the path flow as soon as we encounter a link on the path at which not all path flows have arrived. If some path flows are waiting at link $l$ for other path flows before being propagated, we say that link $l$ is "blocking". As a result, the path flow is propagated along the path until it reaches either its destination or a blocking link.

In our implementation, we have initialized to $-1$ all cumulative arrivals and departures on link $l$ associated with path $p$, denoted by $A_l^p(t)$ and $D_l^p(t)$ respectively, for all links $l$, all paths $p$ that use link $l$, and all time periods $t$. Therefore, the flow can be propagated on link $l$ if all but one path flows have arrived. A procedure for counting the number of path flows that have not arrived to link $l$ at time $t$ is outlined in Procedure 2.

---
**Algorithm 2** Number of path flows that have not arrived to link $(l, t)$

$Count = 0$
**for all** $p \in P : l \in p$ **do**
  **if** path $p$ does not start with $l$ **then**
    $u =$ link preceding $l$ on path $p$
    **if** $D_u^p(t) = -1$ **then**
      $Count = Count + 1$
    **end if**
  **else**
    **if** $A_l^p(t) = -1$ **then**
      $Count = Count + 1$
    **end if**
  **end if**
**end for**
return $Count$

---

Once the path flow has been propagated along link $l$ at time $t$, all links downstream of $l$ are examined, one at a time up to period $t + \tau_l(L_l, t)$, and the flow is propagated

149

along these links if the links are not blocking. In Procedure 3, we maintain a list of active links, $U$, from which we branch out in the network. Once the flow is propagated on some link $l$, we add to the list of active links all links downstream of $l$. Notice that we add several copies of each of them, since the flow that has arrived on link $l$ at time $t$ exits the link some time between the last departure time, $t - 1 + \lceil \tau_l(L_l, t - 1) \rceil$, and the next one, $t + \lfloor \tau_l(L_l, t) \rfloor$.

The procedure described above is repeated for each departure time of each path. We define a list of active paths, $Q$, and an order on the paths to propagate the flow as far as possible into the network. For each path $p$, we denote by $t_p$ the maximum exit time of some path flow on a link in the network. This measures the progression of the flow into the network. A path flow that is blocked in the upstream links of the network will be associated with a small $t_p$, and conversely, a path flow that encounters no blocking link will have a large $t_p$. In Procedure 3, this quantity is updated in the Link Flow Propagation module. Therefore, to avoid blocking links, we choose to process the paths in increasing order of $t_p$. Notice that this order changes dynamically. Once a path has been examined $T$ times, no flow will enter the network, and the path is deleted from the set $Q$.

In Procedure 3, we compute the link travel time for several path flows. So far, we have analyzed the travel time of vehicles with the same destination. In particular, the travel-time function that we derived depends on the ratio of the incoming flow over the exit capacity (see (3.17)). When there are several downstream links, it is not clear how to compute the exit capacity, as it depends on the different priorities (e.g., special lane for right turns). In our simulation, we assumed that there was only a single lane on each arc and no priorities. Accordingly, in our setting, if a flow associated with a particular path is jammed, all flows associated with the other paths are also jammed.

Consider all path flows on link $l$ whose next link is the downstream link $d$. If there were only these path flows on link $l$, the travel time would be a function of the ratio of the sum of these path flows over the capacity of link $d$. If link $d$ is the tightest bottleneck at the exit of link $l$, all other path flows will have the same travel time.

---
**Algorithm 3** Path flow propagation
---
$Q = P$ (list of active paths)

**for all** $p \in P$ **do**

    $A_l^p(t) = D_l^p(t) = -1 \ \forall l, t$

    $t_p = 0$

**end for**

**while** $Q \neq \emptyset$ **do**

    $p = \arg\min_{p \in Q}\{t_p\}$

    $l=$ first link on path $p$

    List $U \leftarrow l$ (list of active links)

    **while** $U \neq \emptyset$ **do**

        $l=$first link in list $U$

        $U \leftarrow U \setminus l$

        $t=$ smallest time period for which $A_l^p(t) = -1$ for some $p \in P$

        **if** (path $p$ starts with $l$) **and** $(t = T)$ **then**

            $Q \leftarrow Q \setminus p$

        **end if**

        **if** Number of path flows that have not arrived to link $(l, t)$=1 **then**

            Travel time computation $(l, t)$

            Link flow propagation $(l, t)$

            **for all** $s = \lceil \tau_l(L_l, t - 1) - 1 \rceil$ to $\lfloor \tau_l(L_l, t) \rfloor$ **do**

                **for all** $d \in D(l)$ **do**

                    $U \leftarrow U \cup d$ (Branching out)

                **end for**

            **end for**

        **end if**

    **end while**

**end while**
---

Therefore, having clustered the path flows by downstream link, the travel time on link $l$ is the maximum of the travel times of all clusters. Kuwahara and Akamatsu (2001) proposed a similar procedure, by allocating the capacity of a downstream link $d$ in proportion of the incoming flows whose next link is link $d$.

In a diverging intersection (one-to-many), the capacity of a downstream link $d$ in period $t + \theta$ is computed as the minimum between the static link capacity $Q_d$ and the dynamic link capacity, if there is queue spillback, $D_d(L_d, t + \theta + 1) - A_d(0, t + \theta) + k_d(L_d, t + \theta)L_d$. In a more general intersection, link $d$ can have several upstream links, other than link $l$. Therefore, the downstream link capacity needs to be reduced to take into account of the flow coming from these other upstream links. If these

incoming flows are not known in time $t$ (because these path flows have not arrived yet to link $d$, being blocked upstream), we estimate these flows from the past flows, taking a linear interpolation.

In the case of a triangular fundamental diagram and with $\theta = \tau_l(L_l, t-1) - 1$, the dynamic effects at the entrance $B(0, t)$ have no influence and the dynamic effects at the road exit $B(L, t)$ are negligible. Therefore, the travel-time function can be fairly accurately approximated as the maximum between the free flow travel time and the travel time with shock (3.17). In fact, if the downstream link $d$ has a bottleneck capacity $CAP_d$ and that a flow $F_d$ is incoming to link $l$ with destination $d$, one can show that there is a shock if and only if $CAP_d(\theta - L_l u_{0,l}) < F_d$, i.e., if $L_l u_{0,l} < \theta + F_d/CAP_d$.

Procedure 4 outlines the computation of the travel time on link $l$ at time $t$. We denote by $D(l)$ and $U(l)$ the set of links downstream of $l$ and upstream of $l$, respectively.

---

**Algorithm 4** Travel time computation $(l, t)$

---

$\theta = \lfloor \tau(L_l, t-1) - 1 \rfloor$
for all $d \in D(l)$ do
    $CAP_d = \min\{Q_d, D_d(L_d, t + \theta + 1) - A_d(0, t + \theta) + k_d(L_d, t + \theta)L_d\}$
    $CAP_d = CAP_d - \sum_{u \in U(d), u \neq l} f_{u,d}(L, t + \theta)$ (remove flows outgoing from other links)
    $F_d = 0$
    for all $p \in P : l, d \in P$ do
        if $p$ starts with $l$ then
            $F_d = F_d + f_p(t)$ (add flow at the path origin)
            $A_l^p(t) = A_l^p(t-1) + f_p(t)$
        else
            $u$=link preceding $l$ on path $p$
            $F_d = F_d + D_u^p(t) - D_u^p(t-1)$ (add flow coming from upstream links)
            $A_l^p(t) = D_u^p(t)$
        end if
    end for
    $\tau_d = \tau_l(L_l, t-1) - 1 + F_d/CAP_d$ (travel time associated with downstream link $d$)
end for
$\tau_l(L_l, t) = \max\{L_l u_{0,l}, \max_{d \in D(l)}\{\tau_d\}\}$

---

Once link travel time have been computed, the flow can be propagated along the

link, following the general idea that $A(0,t) = D(L, t + \tau(L,t))$, detailed in Procedure 5. In addition, Procedure 5 measures the progression of each path flow into the network, $t_p$.

---

**Algorithm 5** Link flow propagation $(l, t)$

---

  **for all** $p \in P : l \in p$ **do**
    **for all** $s = \lceil \tau_l(L_l, t-1) - 1 \rceil$ to $\lfloor \tau_l(L_l, t) \rfloor$ **do**
      $\lambda = \frac{(t+s)-(t-1+\tau_l(L_l,t-1))}{(t+\tau_l(L_l,t))-(t-1+\tau_l(L_l,t-1))}$
      $D_l^p(L_l, t+s) = \lambda A_l^p(0,t) + (1-\lambda) A_l^p(0, t-1)$
    **end for**
    $t_p = \max\{t_p, t + \lfloor \tau_l(L_l, t) \rfloor\}$
  **end for**

---

**Example.** To illustrate the procedure we just described, let us consider the two-link network displayed in Figure 3-10. We denote by link 1 the upstream arc and by link 2 the downstream arc. Link 1 has infinite capacity while link 2 has a capacity of 2. The free flow travel times are equal to 1 for both links. We consider two paths, during two periods of time. Path 1 goes along link 1 and link 2 and has 0 units of flow in the first period, and 2 units of flow in the second period. Path 2 goes along link 2 only and has 0 units of flow in the first period and 1 unit of flow in the second period.



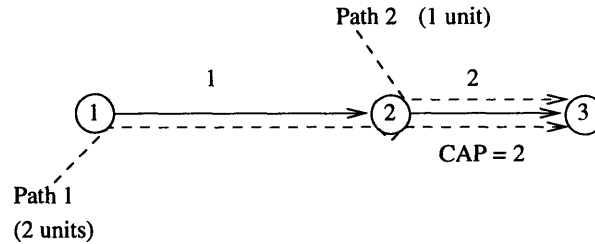Figure 3-10: A two-link network example.

We initialize $Q = \{1, 2\}$. In the first iteration, the two measures of progression of the paths, $t_1$ and $t_2$ are set to zero. We arbitrarily select path 1 to start with. In period 1, path 1 has zero units of flow that will travel at the free flow speed; therefore, $A_1^1(1) = 0 = D_1^1(s)$, for $s \leq 2$, and $t_1$ is updated to 2. The flow cannot be propagated

on link 2 because the flow from path 2 has not arrived yet. We say that link 2 is blocking.

Since $0 = t_2 < t_1 = 2$, we select path 2 in the second iteration. Path 2 has zero units of flow in the first period, $A_2^2(1) = 0$. The flow on link 2 in period 1 is the sum of flows of paths 1 and 2, i.e., $F = f_2(1) + (D_1^1(1) - D_1^1(0)) = 0$. This zero flow will travel at the free flow speed. Therefore, $D_2^p(s) = 0$, for $p = 1, 2$ and $s \leq 2$; $t_2$ is updated to 2.

Since $t_2 = t_1 = 2$, we pick arbitrarily path 2 in the third iteration. The flow on link 2 in period 2 is the sum of flows of paths 1 and 2, i.e., $F = f_2(2) + (D_1^1(2) - D_1^1(1)) = 1 + 0 = 1$ and will travel at the free flow speed. Therefore, $D_2^1(3) = 0$, and $D_2^2(3) = 1$. As path 2 has been examined twice, it is deleted from $Q$.

The set of active paths $Q$ contains only one element, path 1. The incoming flow is 2, and $A_1^1(2) = 2$. The downstream capacity is equal to 2, from which we subtract the flow that came from path 2 in period $t - 1 + \tau_1(L_1, t - 1) = 2$, equal to 1. As a result, only 1 unit of capacity is available for the flow of path 1. The travel time is the maximum between the free flow travel time, 1, and the travel time with shock, $\tau_1(L_1, 1) - 1 + 2/1 = 2$. Therefore, $D_1^1(4) = 2$ and $D_1^1(3) = 1$. The list of active links contains two copies of link 2, corresponding to the propagation of flow in periods 3 and 4.

Since no flow comes from path 2 in period 3, the flow from link 1 can be propagated downstream. As a result, $A_2^1(3) = 1$. Since the flow propagates at the free flow speed on the second link, $D_2^1(4) = 1$.

Likewise, since no flow comes from path 2 in period 4, the flow from link 1 can be propagated downstream. As a result, $A_2^1(4) = 2$. Since the flow propagates at the free flow speed on the second link, $D_2^1(5) = 2$.

### 3.5.3 Numerical Example

To validate the procedure we introduced for solving DUE problems, we consider the transportation network test example proposed by Kuwahara and Akamatsu (2001) with time-varying demands. As shown in Figure 3-11, this transportation network of 6

nodes and 7 links is a simple model of a freeway and an arterial running parallel. The link data are shown in Table 3.1, associated with a triangular fundamental diagram. We consider the evening commute, where users are going back from their work place to the residential area, by taking either the freeway or the arterial. In particular, we consider two origin-destination pairs, $\{1,2\}$ and $\{1,3\}$ with the following traffic demand:

$$d^{\{1,2\}}(t) = \begin{cases} 1000 \text{ veh/hour} & 0 \leq t \leq 1 \text{ (h) and } 3 < t \leq 5 \text{ (h)} \\ 2000 \text{ veh/hour} & 1 < t \leq 3 \text{ (h)} \end{cases}$$

$$d^{\{1,3\}}(t) = \begin{cases} 2000 \text{ veh/hour} & 0 \leq t \leq 1 \text{ (h) and } 3 < t \leq 5 \text{ (h)} \\ 4000 \text{ veh/hour} & 1 < t \leq 3 \text{ (h)}. \end{cases}$$



Figure 3-11: An example network

Table 3.1: Link characteristics

| Link | $u_0$ (h/km) | $w_0$ (h/km) | $k^{max}$ (veh/km) | Length (km) | $f^{max}$ (veh/hour) |
|---|---|---|---|---|---|
| (1,4) | 0.025 | 0.05 | 450 | 2 | 6000 |
| (4,5) | 0.0125 | 0.05 | 375 | 16 | 6000 |
| (5,6) | 0.0125 | 0.05 | 250 | 8 | 4000 |
| (6,3) | 0.025 | 0.05 | 225 | 2 | 3000 |
| (1,2) | 0.025 | 0.05 | 450 | 16 | 6000 |
| (2,3) | 0.025 | 0.05 | 450 | 12 | 6000 |
| (5,2) | 0.025 | 0.05 | 450 | 2 | 6000 |

In free flow, the freeway is the most attractive path to go from the work place to the residential area. For instance, for the OD pair $\{1,3\}$, the freeway free-flow travel time is 0.4 h while the arterial free flow travel time is 0.7 h. On the other

hand, in the peak traffic period, link $(6, 3)$ will be saturated, since it has a capacity of 3000 vehicles/hour, while the peak demand is 4000 vehicles/hour. As a result, a queue will appear and propagate backwards, and the freeway will become less and less attractive, in comparison to the arterial. As shown in Kuwahara and Akamatsu (2001), the impact of the queue spillback is so significant that a point queue model would not be appropriate for this example.

In our simulations, we considered 5 hours of traffic demand, discretized in time steps of 0.01 h, that is, we considered 500 time periods.

To solve the dynamic user equilibrium problem on this network, we implemented in C++ a dynamic version of the Frank-Wolfe algorithm. At every iteration we solved a linear optimization problem, calling the CPLEX library, to solve the path formulation of the DUE, described in Subsection 5.1. The Frank-Wolfe algorithm for solving static variational inequality problems is fairly standard in the literature (e.g., see Martos 1975), and we outline it only briefly in Algorithm 6.

---

**Algorithm 6** A dynamic Frank-Wolfe algorithm for solving (3.20)

Choose $\epsilon > 0$
$m = 0$,
$\mathbf{f}_p^m = 0 \forall p \in P$
$\tau^m = \tau(\mathbf{f}^m)$ (free flow travel times)
**for all** $w \in W$ **do**
  $p =$ shortest path w.r.t. $\tau^m$
  $f_p^m(t) = d^w(t), \forall t = 1, ..., T$
**end for**
**repeat**
  $m \leftarrow m + 1$
  $\tau^m = \tau(\mathbf{f}^m)$ (see Section 3.5.2)
  $\mathbf{f}^* = \arg\max_{\mathbf{f}} \{(\tau^m)'\mathbf{f} : \sum_{p \in P^w} f_p^w(t) = d^w(t), \forall w \in W, \mathbf{f} \geq 0\}$ (Linear Optimization)
  **if** $\tau(\mathbf{f}^*)'(\mathbf{f}^m - \mathbf{f}^*) < 0$ **then**
    Use binary search to find $\mathbf{f}^{m+1} = \lambda\mathbf{f}^m + (1 - \lambda)\mathbf{f}^*$, $0 \leq \lambda \leq 1$, such that $\tau(\mathbf{f}^{m+1})'(\mathbf{f}^m - \mathbf{f}^*) = 0$
  **end if**
**until** $\tau(\mathbf{f}^m)'(\mathbf{f}^m - \mathbf{f}^*) < \epsilon$

---

As shown in Figure 3-12, the algorithm made significant progress in the first iterations (in terms of the decrease of $\tau(\mathbf{f}^m)'(\mathbf{f}^m - \mathbf{f}^*)$) but had trouble to converge to

zero thereafter. In fact, as we mentioned above, the path travel times are not necessarily monotone, and the algorithm is not a priori guaranteed to converge. However, from Figure 3-12, it seems that the algorithm would ultimately converge, but not monotonically.



Figure 3-12: Convergence Plot of the Frank-Wolfe Algorithm

The key issue in implementing this algorithm is the evaluation of the path travel times that relies on the procedure described in Section 3.5.2. Even if the procedure for disaggregating path flows into link flows may seem cumbersome, it runs pretty fast. In fact, the first 30 iterations took less than a minute.

Figures 3-13 and 3-14 display the curves of cumulative number of vehicles at the entrance of each link, for the OD pairs $\{1,3\}$ and $\{1,2\}$ respectively. Since the path choice decisions are based on the actual travel time, and not the instantaneous travel time as in Kuwahara and Akamatsu (2001), the figures are very different from Figures 6 and 7 in their paper.

For OD pair $\{1,3\}$, the freeway is the first used path. As in Kuwahara and Akamatsu, a queue completely covers link $(6,3)$ after 1.35 hour and spills back onto the previous link. After 1.75 hour, the queue also covers link $(5,6)$ and spills back onto link $(4,5)$. Notice that, since the capacity of link $(6,3)$ is fixed and that there is no incoming flow into the freeway at the intermediate nodes, the density of the queue remains constant, and so is the travel time when the queue covers an entire link.

One significant difference with their model is that path $\{(1,2),(2,3)\}$ starts being

157

Figure 3-13: Cumulative link arrivals for the OD pair $\{1, 3\}$

used after 1.9 hour, and not after 2.09 hours. In fact, our model is based on the actual travel time and not the instantaneous travel time, and drivers anticipate earlier the possible delays on the freeway. Another difference with their results is that the two paths are never used simultaneously, while in Figure 3-13, the curves of cumulative arrivals into $(1, 4)$ and $(1, 2)$ keep growing together. Specifically, in their model, there is no flow incoming to link $(1, 4)$ in the time intervals $[2.09, 2.29]$ and $[2.85, 3.15]$, as everything goes through the second path. The sudden switch of used paths that they observed seems to be a consequence of the path choice based on the instantaneous travel times. In their model, even though everybody decides to take an alternate path, the instantaneous travel time of the original path keeps increasing, deterring future flow to choose this path.

For OD pair $\{1, 2\}$ (see Figure 3-14), the effects of the queue originating from link $(6, 3)$ have an impact on the travel time of link $(4, 5)$ in period 1.75, as the queue backs up into link $(4, 5)$. Because the travel time is increasing, the arterial is becoming more and more attractive. In period 1.79, the arterial is so attractive that all flow is taking it. In contrast, in Kuwahara and Akamatsu's model, the drivers did not anticipate the queue on the freeway, and the switch to the arterial only occurs after 2 hours. Moreover, in our model, the arterial is being used during the whole peak

158

Figure 3-14: Cumulative link arrivals for the OD pair $\{1, 2\}$

period, as the queue on the freeway does not decrease. With the instantaneous travel time decisions, since there are two switches from the freeway to the arterial, and then from the arterial to the freeway for OD pair $\{1,3\}$ during the peak period, the travel time on the freeway decreases after 2.34 hours, making the arterial less attractive.

In addition, we can observe that the principle of optimality is maintained dynamically. For example, if one of the shortest paths for OD $\{1,3\}$ in period 2 is $\{(1,2),(2,3)\}$, one of the shortest paths for OD $\{1,2\}$ in the same period is $\{(1,2)\}$. This observation may be used to speed up the solution of the DUE problem.

## 3.6 Conclusions

In this chapter, we proposed a methodology for deriving an analytical travel-time function based on the theory of kinematic waves. In particular, we derived a travel-time function in the cases of a triangular and a quadratic fundamental diagrams. The derived travel-time function integrates first-order traffic dynamics as well as shock waves. Moreover, we illustrated numerically that this travel-time function is very consistent with the simulated travel times proposed in the literature.

With the advent of advanced transportation management systems (ATMS), a lot

159

of attention has been devoted to solving the dynamic user equilibrium problem. We showed that the travel-time function that we derived can be incorporated within a DUE setting, since it is (strictly) monotone and satisfies the FIFO property under certain conditions. As an illustration, we applied our model to a simple evening commute problem, emphasizing the dynamic nature of the traffic assignment.

Further research is necessary to develop a mathematical formulation of the flow propagation proposed in Section 3.5. A mathematical formulation of the problem, defined in terms of links instead of paths, would take advantage of the continuity and the (strict) monotonicity of the travel-time functions derived in Section 3.4. Moreover, the model needs to be enriched to capture local triggers of congestion such as deadlocks (on cyclic networks), traffic lights, multi-lane roads, etc. Finally, the proposed model for delays can be incorporated into more general or alternative traffic assignment problems as the system equilibrium, the stochastic DUE, or the DUE with departure time choices.

# Conclusions

This thesis investigates the impact of lack of information and decentralization of decision-making on the efficiency of inventory, supply chain, and transportation systems. The first part of the thesis quantifies the loss of profit in the newsvendor model with limited knowledge about the demand probability distribution (Chapter 1) and with several competing decision-makers (Chapter 2). The second part of the thesis analyzes the dynamic traffic equilibrium in a transportation network, when drivers choose their itinerary selfishly (Chapter 3).

It is common to hear that "all models are wrong." One of the greatest challenges for modelers is to trade off the complexity of reality with mathematical tractability. On the one hand, a model that fails to capture the main constituents of a problem will lead to poor decision-making. On the other hand, a mathematical problem that is impossible to solve is useless. Both the newsvendor problem and the traffic equilibrium problem are traditionally associated with a certain set of assumptions (e.g., probabilistic description of demand and centralized decision-making in the newsvendor model; static representation of the traffic equilibrium problem). One of the goals of this thesis was to analyze the performance of these models under more general assumptions.

Our motivation for relaxing assumptions in these two models was threefold. First, a model that relies on fewer assumptions is usually a more accurate representation of reality and can therefore lead to better decisions. For instance, in transportation, congestion is inherently dynamic; therefore, a traffic equilibrium model that captures the dynamic interaction of traffic flows is more informative than a static model, which would focus only on the peak congestion period. In particular, in a dynamic model,

travelers anticipate congestion ahead and choose their itinerary accordingly. Similarly, in supply chain management, decision-making is decentralized; incorporating the multiplicity of decision-makers into an inventory model sheds light on the inefficiency of simple contracting mechanisms. In this study, we analyzed various supply chain configurations to develop insight into the factors that accentuate or mitigate this loss of efficiency.

Our second goal was to test the robustness of the classic models under more general assumptions. For instance, the newsvendor model assumes that the demand probability distribution is known but says nothing about how to estimate it. By assuming only limited information about the demand, we have derived guidelines for selecting a demand distribution, as a function of the prior information about the demand. In particular, the uniform demand distribution should be chosen when only the range is known and the exponential distribution should be chosen when only the mean is known. When both the mean and the variance are known, the normal distribution can be assumed only if the coefficient of variation is less than .3; for larger coefficients of variation, the gamma distribution seems to be robust.

Finally, and perhaps paradoxically, making fewer assumptions sometimes simplifies the mathematical analysis. For instance, the multi-item newsvendor problem with capacitated constraints is a difficult nonlinear problem, often prone to numerical instabilities. On the other hand, when only the ranges of the demand distributions are specified, the problem can be formulated as a simple linear optimization problem, despite the fact that it still captures the stochastic nature of the demand. More generally, this thesis illustrates how adopting a new perspective on a specific problem might lead to more relevance, more insights, and more efficient solution techniques.

162

# Appendix A

# Proofs

## A.1 Chapter 1

### Proof of Proposition 1

(a) A function $f$ defined on a convex set $C$ is quasi-concave if $f(\lambda x_1 + (1 - \lambda)x_2) \geq$ $\min\{f(x_1), f(x_2)\}$ for $x_1, x_2 \in C$ (Boyd and Vandenberghe 2004).

Let us consider the case $y \leq z$. We will show that the function $\Phi(z; y)$ is quasi-concave for $z \in [y, \infty)$. Fix two points $z_2 > z_1 > y$, and a scalar $\lambda \in (0, 1)$, and let $z_0 = \lambda z_1 + (1 - \lambda)z_2$. The optimal solution of problem (1.4) when $z = z_i$, is denoted by $F_i \in \mathcal{D}$, for $i = 0, 1, 2$.

The objective function of (1.4) can be rewritten as $\int_y^z \bar{F}(x)dx$. Therefore, $F_0(z_0) \geq F_1(z_1)$ by optimality of $F_1$, and $F_0(z_0) \leq F_2(z_2)$ by optimality of $F_0$. Consequently, $F_2(z_2) \geq F_0(z_0) \geq F_1(z_1)$.

If $F_1(z_1) > 1 - c/p$, $F_2(z_2) \geq F_0(z_0) > 1 - c/p$. As a result, the quasi-concavity of $\Phi(z; y)$ is established as follows:

$$
\begin{aligned}
\Phi(z_0; y) - \Phi(z_2; y) &= p \int_\Omega (\min\{x, z_0\} - \min\{x, y\})dF_0(x) - c(z_0 - y) \\
&\quad -p \int_\Omega (\min\{x, z_2\} - \min\{x, y\})dF_2(x) + c(z_2 - y), \\
&\geq -p \int_{z_0}^{z_2} \bar{F}_2(x)dx + c(z_2 - z_0), \\
&\geq -p(z_2 - z_0)\bar{F}_2(z_2) + c(z_2 - z_0),
\end{aligned}
$$

163

$$\geq \quad -p(z_2 - z_0)(c/p) + c(z_2 - z_0) > 0,$$

where the first inequality follows from optimality of $F_0$ when $z = z_0$, the second inequality follows from the fact that $F_2$ is increasing, and the third inequality comes from the assumption $F_2(z_2) > 1 - c/p$.

Similarly, if $F_2(z_2) < 1 - c/p$, $F_1(z_1) \leq F_0(z_0) < 1 - c/p$. Then, for similar reasons as above,

$$
\begin{aligned}
\Phi(z_0; y) - \Phi(z_1; y) &\geq \quad p \int_{z_1}^{z_0} \bar{F}_1(x)dx - c(z_0 - z_1), \\
&\geq \quad p(z_0 - z_1)\bar{F}_1(z_0) - c(z_0 - z_1), \\
&\geq \quad p(z_0 - z_1)(c/p) - c(z_0 - z_1) > 0.
\end{aligned}
$$

Finally, when $F_1(z_1) \leq 1 - c/p$ and $F_2(z_2) \geq 1 - c/p$, $\Phi(z_0; y) \geq \Phi(z_1; y)$ if $F_0(z_0) \leq 1 - c/p$, and $\Phi(z_0; y) \geq \Phi(z_2; y)$ if $F_0(z_0) \geq 1 - c/p$.

Hence $\Phi(z; y)$ is quasi-concave for $z \geq y$. One can apply the same argument to show that the function is also quasi-concave for $z \leq y$, but we omit the details for the sake of brevity.

However, $\Phi(z, y)$ is not necessarily quasi-concave on $[0, \infty)$. For instance, consider $\mathcal{D}$ as the set of all distributions with support $\Omega$. Then, $\Phi(z; y)$ equals $(p - c)(z - y)$ when $z \geq y$, and $-c(z - y)$ otherwise. As a result, $\Phi(z; y)$ is piecewise linear convex, with a kink at $z = y$.

(b) For a fixed demand distribution $D \in \mathcal{D}$, $E[\Pi(z, D) - \Pi(y, D)]$ is a convex function of $y$, since $-E[\Pi(y, D)]$ is the negative of the objective of the classic newsvendor problem (1). Because convexity is preserved under maximization (Porteus 2002), the maximum of the regret over all demand distributions is a convex function of $y$. $\quad \square$

## Proof of Theorem 1

Consider problem (1.4) with only the normalization constraint, i.e., $\int_l^u dF(x) = 1$. From Proposition 1, two cases need to be considered: when $y \leq z$ and when $y \geq z$.

When $y \leq z$, the worst-case demand distribution that solves (1.4) is a unit impulse at $z$. Accordingly, the optimal value of (1.4) equals $(p - c)(z - y)$. Maximizing the value of the regret over all feasible $z \in [y, u]$, we obtain a regret equal to $(p-c)(u-y)$.

Similarly, when $y \geq z$, the worst-case demand distribution is also a unit impulse at $z$, and gives rise to a regret of $-c(z - y)$. The maximum regret, taken over all feasible $z \in [l, y]$, equals $-c(l - y)$.

From Proposition 1, the optimal order quantity $y$ equates the two maximum regrets, leading to the theorem statement. $\qquad \square$

## Proof of Theorem 2

Problem (1.4) can be formulated as a semi-infinite linear optimization problem. By strong duality (Smith 1995), the primal problem is equivalent to the following dual problem:

$$\min_{\alpha_0, \alpha_1} \quad \alpha_0 + \alpha_1 \mu, \tag{A.1}$$
$$\text{s.t.} \quad \alpha_0 + \alpha_1 x \geq \min\{x, z\} - \min\{x, y\}, \quad \forall l \leq x \leq u.$$

(a) If $z \geq y$, a dual feasible function is any straight line, possibly discontinuous at $z$, with ordinate $\alpha_0$ and slope $\alpha_1$ that is nonnegative for all $x \geq l$, lies above the line $x - y$ between $y$ and $z$, and above the line $z - y$ for all $z \leq x \leq u$. In an optimal solution, the constraints are tight at $z$. (Otherwise, one could reduce the cost, by lowering either $\alpha_0$ or $\alpha_1$, without violating a constraint.) It is easy to see that there are two possible optimal solutions: either the straight line that goes through the points $(l, 0)$ and $(z, z - y)$, or the horizontal line at $z - y$. In the first case, $\alpha_0 = -l\alpha_1$, $\alpha_1 = (z - y)/(z - l)$, and the optimal value of (A.1) is equal to $(z - y)(\mu - l)/(z - l)$. In the second case, $\alpha_0 = z - y$, $\alpha_1 = 0$, and the optimal value of (A.1) is equal to $(z - y)$. Therefore, the first case is optimal if and only if $\mu \leq z$.

(a.1) If $\mu \leq z$, the regret is equal to $(z - y)(p(\mu - l)/(z - l) - c)$, which is concave in $z$. The regret, optimized over all possible values of $z \in [\mu, u]$, is maximized at $z^* = l + \sqrt{(y - l)(\mu - l)(p/c)}$, if $\mu \leq z^* \leq u$, at $\mu$ if $\mu \geq z^*$, and at $u$ if $z^* \geq u$.

Substituting $z$ by its optimal value simplifies the regret to

$$\begin{cases} (\mu - y)(p - c), & \text{if } y \leq l + \frac{c}{p}(\mu - l), \\ (l - y + \sqrt{(y - l)(\mu - l)\frac{p}{c}})(\sqrt{cp\frac{\mu - l}{y - l}} - c), & \text{if } l + \frac{c}{p}(\mu - l) \leq y \leq l + \frac{c}{p}\frac{(u-l)^2}{\mu - l}, \\ (u - y)(p\frac{\mu - l}{u - l} - c), & \text{if } y \geq l + \frac{c}{p}\frac{(u-l)^2}{\mu - l}. \end{cases}$$

(a.2) If $z \leq \mu$, the regret is equal to $(z - y)(p - c)$. The maximum regret, when it is optimized over $z \in [l, \mu]$, is attained at $z = \mu$ and equal to $(\mu - y)(p - c)$.

(b) On the other hand, if $y \geq z$, the right hand side of the constraints is nonincreasing. It is easy to see that the optimal solution is $\alpha_0 = \alpha_1 = \alpha_2 = 0$. Therefore, the regret equals $-c(z - y)$. The maximum regret, optimized over all values of $z \in [l, y]$, is equal to $-c(l - y)$.

From Proposition 1, the quantity $y$ balances the opportunity cost from ordering too much with the opportunity cost from ordering too little. If $y \leq l + (c/p)(\mu - l)$, the optimal order quantity minimizes the maximum of the two following convex functions:

$$\min_{y \geq 0} \max\{c(y - l), (p - c)(\mu - y)\}.$$

In this case, it is optimal to order $y^* = (c/p)l + ((p - c)/p)\mu$. The condition $y^* \leq l + (c/p)(\mu - l)$ reduces to $(p/c) \leq 2$.

If $l + (c/p)(\mu - l) \leq y \leq l + (c/p)(u - l)^2/(\mu - l)$, the order quantity minimizes the following expression:

$$\min_{y \geq 0} \max\{c(y - l), (l - y + \sqrt{(y - l)(\mu - l)\frac{p}{c}})(\sqrt{cp\frac{\mu - l}{y - l}} - c)\}.$$

The minimum is at the intersection of the two curves, i.e., when $y^* = l + (p/c)(\mu - l)/4$. For this value of $y$, the conditions $l + (c/p)(\mu - l) \leq y \leq l + (c/p)(u - l)^2/(\mu - l)$ translate to $2 \leq (p/c) \leq 2(u - l)/(\mu - l)$.

Finally, if $y \geq l + (c/p)(u - l)^2/(\mu - l)$, the optimal order quantity solves the

following:

$$\min_{y \geq 0} \max\{c(y - l), (u - y)(p\frac{\mu - l}{u - l} - c)\}.$$

The minimum occurs at the intersection of the two curves, i.e., when $y^* = u - \frac{c}{p}\frac{(u-l)^2}{\mu-l}$. For this value of $y$, the condition $y \geq l + (c/p)(u - l)^2/(\mu - l)$ translates to $(p/c) \geq 2(u - l)/(\mu - l)$.

Finally, the Price of Information is obtained by plugging the optimal values of y into the regret functions. □

## Proof of Theorem 3

Problem 1.4 can be formulated as the following semi-infinite linear optimization problem:

$$\begin{aligned}
\max_{P(x)} \quad & \int_0^\infty (\min\{x, z\} - \min\{x, y\})P(x)dx, \\
\text{s.t.} \quad & \int_0^\infty P(x)dx = 1, \\
& \int_0^\infty xP(x)dx = \mu, \\
& \int_\mu^\infty P(x)dx = \tfrac{1}{2}, \\
& P(x) \geq 0,
\end{aligned}$$

where the next to last constraint ensures that the mean is the median of the distribution. Rather than solving the primal problem, we solve its dual problem (1.5) geometrically. We denote by $\alpha_i$ the dual variable associated with the $i$th constraint.

Since the mean is equal to the median, then $z \leq \mu$ if and only if $c/p \geq 1/2$. As a result, it is optimal to order $y \leq \mu$ if and only if $c/p \geq 1/2$. Therefore, four cases need to be considered, depending on whether $y \geq z$ or not, and whether $c/p \geq 1/2$ or not.

**Case 1: $c/p \geq 1/2$.** If $z \geq y$, a dual feasible solution is a straight line with ordinate $\alpha_0$, slope $\alpha_1$, and possibly discontinuous at $\mu$ by an amount $\alpha_2$, that is nonnegative for all $x \geq 0$, lies above the line $x - y$ between $y$ and $z$, and above the line $z - y$ for all

167

$z \leq x$. The optimal dual solution is a horizontal line, with $\alpha_0 = z - y$, $\alpha_1 = \alpha_2 = 0$. The associated regret, $(p - c)(z - y)$, is maximized at $z = \mu$.

If $z \leq y$, the dual constraints are piecewise linear decreasing. The optimal dual solution is a horizontal line at zero, discontinuous at $\mu$ by an amount $\alpha_2 = z - y$. The regret is equal to $(p/2 - c)(z - y)$ and maximized when $z = 0$.

The robust order quantity equates the following maximum regrets:

$$\min_{y \geq 0} \left\{ (p - c)(\mu - y), (p/2 - c)(-y) \right\},$$

and is equal to $y = 2\mu(p - c)/p$.

**Case 2:** $c/p \leq 1/2$. The dual feasible sets are the same as when $c/p \geq 1/2$. If $z \geq y$, there are two possible optimal solutions, depending on whether $z \geq 2\mu$ or not. If $z \leq 2\mu$, the optimal dual solution is also a straight line, but with $\alpha_2 = z - y$ and $\alpha_0 = \alpha_1$. The associated regret equals $(p/2 - c)(z - y)$ and is maximized when $z = 2\mu$. On the other hand, when $z \geq 2\mu$, the optimal dual solution is a piecewise linear line, with slope $\alpha_1 = (z - y)/(z - \mu)$, discontinuous at $\mu$ by an amount $\alpha_2 = -\alpha_1\mu$. The associated regret equals $(p\mu/(2z - 2\mu) - c)(z - y)$ and is maximized when $z = \mu + 1/(2c)\sqrt{(2cp\mu)(y - \mu)}$. The maximizing $z$ is greater than $2\mu$ whenever $y \geq \mu(p + 2c)/p$.

On the other hand, when $z \leq y$, the optimal dual solution is a horizontal straight line at zero. The related regret $-c(z - y)$ is maximized when $z = \mu$.

If $y \leq \mu(p + 2c)/p$, the robust order quantity equates the following maximum regrets:

$$\min_{y \geq 0} \left\{ (p/2 - c)(2\mu - y), -c(\mu - y) \right\},$$

and is also equal to $y = 2\mu(p - c)/p$. Condition $y \leq \mu(p + 2c)/p$ then simplifies to $p \leq 4c$.

If on the other hand $y \geq \mu(p + 2c)/p$, the robust order quantity minimizes the

following regrets:

$$\min_{y \geq 0} \left\{ (\frac{p\mu}{2(z - \mu)} - c)(z - y), -c(\mu - y) \right\},$$

with $z = \mu + 1/(2c)\sqrt{(2cp\mu)(y - \mu)}$, and is equal to $\mu(p + 8c)/(8c)$. Condition $y \geq \mu(p + 2c)/p$ simplifies then to $p \geq 4c$. $\qquad\square$

## Proof of Theorem 4

Following Popescu (2005), the closed convex set of symmetric distributions $\mathcal{D}$ can be generated by pairs of symmetric Diracs. Using this characterization, the dual problem (1.5) can be formulated as (Popescu 2005):

$\min_{\alpha_0, \alpha_1} \quad \alpha_0 + \alpha_1 \mu,$

s.t. $\quad 2\alpha_0 + 2\mu\alpha_1 \geq$

$$\min\{\mu - x, z\} + \min\{\mu + x, z\} - \min\{\mu - x, y\} - \min\{\mu + x, y\},$$

$$\forall 0 \leq x \leq \mu.$$

The dual problem can easily be solved geometrically. A dual feasible solution is a horizontal line, lying above the piecewise linear function described by the right-hand side of the constraint. By symmetry, the mean is equal to the median. As a result, $y, z \geq \mu$ if and only if $c/p \leq 1/2$. Four cases need to be considered, depending on whether $z \geq y$ or not, and $c/p \geq 1/2$ or not.

**Case 1:** $c/p \leq 1/2$. When $z \geq y$, the ordinate of the line is equal to $z - y$. The associated regret, $(p/2 - c)(z - y)$, is maximized at $z = 2\mu$.

When $z \leq y$, the ordinate of the line is equal to zero. The associated regret, $-c(z - y)$, is maximized at $z = \mu$.

Equating both regrets, $(p/2 - c)(2\mu - y) = -c(\mu - y)$, gives rise to the robust order quantity $y = 2\mu(p - c)/p$.

**Case 2:** $c/p \geq 1/2$. When $z \geq y$, the ordinate of the line is equal to $2(z - y)$. The associated regret, $(p - c)(z - y)$, is maximized at $z = \mu$.

When $z \leq y$, the ordinate of the line is equal to $z - y$. The associated regret, $(p/2 - c)(z - y)$, is maximized when $z$ is zero.

Equating both regrets, $(p - c)(\mu - y) = -(p/2 - c)y$, gives rise to the robust order quantity $y = 2\mu(p - c)/p$. $\square$

## Proof of Theorem 5

Following Popescu (2005), the closed convex set of unimodal distributions with mode $m$, $\mathcal{D}$, can be generated with $m$-rectangular distributions (i.e., uniform distributions over a segment bounded by $m$). With this representation, the dual problem (1.5) can be formulated as (Popescu 2005):

$$\min_{\alpha_0} \quad \alpha_0,$$
$$\text{s.t.} \quad \alpha_0(m - x) \geq \int_x^m \min\{\xi, z\} - \min\{\xi, y\}d\xi, \quad \forall 0 \leq x \leq m,$$
$$\alpha_0(x - m) \geq \int_m^x \min\{\xi, z\} - \min\{\xi, y\}d\xi, \quad \forall m \leq x \leq u.$$

The dual problem can easily be solved geometrically. A dual feasible solution is a piecewise linear function, passing through $(m, 0)$, with slope $-\alpha_0$ before $m$ and $\alpha_0$ after $m$, lying above the piecewise quadratic function described by the right-hand side of the constraint. Six cases need to be considered, depending on the relative order of $z$, $y$, and $m$.

**Case 1:** $z \leq y \leq m$. The right-hand side is constant linear for $x \leq z$, quadratic between $z$ and $y$, then increasing linear until $m$, and decreasing linear thereafter. The optimal dual solution is such that the constraint is tight at zero. Therefore, $\alpha_0 = (z - y)(m - z/2 - y/2)/m$, and the regret equals $p(z - y)(m - z/2 - y/2)/m - c(z - y)$. The maximum regret is attained at $z = m(p - c)/p$ or $z = y$, whichever is the smallest. It is equal to zero when $y \leq m(p - c)/p$, and to $(yp - m(p - c))^2/(2pm)$ otherwise.

170

**Case 2:** $z \leq m \leq y$. The right-hand side is constant up to $z$, then quadratic between $z$ and $y$, with a change of concavity at $m$, and linear after $y$. An optimal dual solution is such that the constraint is tight at zero. Therefore, $\alpha_0 = -1/2(m-z)^2/m$, and the regret equals $-p/2(m-z)^2/m - c(z-y)$. The maximum regret equals $c/(2p)(mc - 2mp + 2yp)$ and is attained at $z = m(p-c)/p$.

**Case 3:** $m \leq z \leq y$. The right-hand side is zero for $x \leq z$ and negative thereafter. The optimal dual solution is equal to zero, and the regret equals $-c(z-y)$. The regret is maximized when $z = m$; as a result, Case 3 is dominated by Case 2.

**Case 4:** $m \leq y \leq z$. The right-hand side is zero for $x \leq y$ and increasing after; it is convex between $y$ and $z$, and linear beyond $z$. The optimal solution is such that the function $\alpha_0(x-m)$ crosses the constraint set at $u$, i.e., $\alpha_0 = 1/(u-m)(z-y)(u-z/2-y/2)$. The regret equals $(p/(u-m)(u-z/2-y/2)-c)(z-y)$ and is maximized at $z = u(p-c)/p+mc/p$. At its maximum, the regret equals $(u(c-p)-mc+yp)^2/(2p(u-m))$.

**Case 5:** $y \leq m \leq z$. The right-hand side is constant for $x \leq y$, concave decreasing between $y$ and $m$, convex increasing between $m$ and $z$ and linear increasing after. The optimal solution is such that the function $\alpha_0(x-m)$ crosses the constraint set at $u$, i.e., $\alpha_0 = 1/(u-m)((z-m)(z/2+m/2-y)+(z-y)(u-z))$. The maximum regret, attained at $z = u(p-c)/p+c/pm$, is equal to $(p-c)((u+m)/2-c/p(u-m)/2-y)$.

**Case 6:** $y \leq z \leq m$. The right-hand side is constant for $x \leq y$, concave decreasing between $y$ and $z$, linear decreasing until $m$, and linear increasing thereafter. The optimal solution is to have the slope of the function equal to the slope of the linear piece of the constraint. Accordingly, $\alpha_0 = z - y$; the regret equals $(p-c)(z-y)$ and is maximized at $z = m$; as a result, Case 6 is dominated by Case 5.

The robust order quantity equates the regrets. Taking $y \leq m(p-c)/p$ is suboptimal, since the regret when $z \geq y$, equal to $(p-c)(u-y)$ dominates the regret when

171

$z \leq y$, equal to zero. When $m(p-c)/p \leq y \leq m$, $y$ minimizes the maximum regrets:

$$\min_{y \geq 0} \max\{\frac{(yp - m(p-c))^2}{2pm}, (p-c)(\frac{u+m}{2} - \frac{c}{p}\frac{u-m}{2} - y)\},$$

and is equal to $\sqrt{m(p-c)(2cm + u(p-c))}/p$. The condition that $y \geq m(p-c)/p$ is automatically satisfied if $u \geq m(1 - 2c/(p-c))$, and the condition that $y \leq m$ simplifies to $u \leq m(1 + (\frac{c}{p-c})^2)$.

When $y \geq m$, the robust order quantity minimizes the maximum regrets:

$$\min_{y \geq 0} \max\{\frac{c}{2p}(mc - 2mp + 2yp), \frac{(u(c-p) - mc + yp)^2}{2p(u-m)}\}$$

and is equal to $u - \sqrt{c(u-m)(2up - cu + 2cm - 2mp)}/p$. The condition that $y \geq m$ simplifies to $u \geq m(1 + (\frac{c}{p-c})^2)$.                     □

## Proof of Theorem 6

Following Popescu (2005), the set of unimodal and symmetric distributions with mean $\mu$ can be generated using a mixture of $\mu$-centered rectangular distributions (i.e., uniform distributions centered around $\mu$). Using this representation, the dual problem (1.5) can be formulated as follows:

$$\min_{\alpha_0, \alpha_1} \quad \alpha_0 + \alpha_1 \mu,$$
$$\text{s.t.} \quad 2t(\alpha_0 + \mu\alpha_1) \geq \int_{\mu-t}^{\mu+t} \min\{\xi, z\} - \min\{\xi, y\}d\xi, \quad \forall 0 \leq t \leq \mu.$$

A dual feasible solution is the slope of a straight line, passing through the origin originating. Because the mean equals the median (by symmetry), $z \geq \mu$ whenever $c/p \leq 1/2$. As a result four cases need to be considered, depending on whether $z \geq y$ or not, and whether $c/p \geq 1/2$ or not.

**Case 1:** $c/p \geq 1/2$. When $z \geq y$, the right-hand side of the dual constraint is increasing, linearly with slope $2(z-y)$ for $t \leq \mu - z$, then concavely between $\mu - z$ and $\mu - y$, and then linearly with slope $z - y$. The dual optimal solution is a straight

172

line with slope equal to the first piece of the constraint, i.e., $\alpha_0 + \alpha_1\mu = z - y$. The regret, $(p - c)(z - y)$, is maximized at $z = \mu$.

If on the other hand $z \leq y$, the right-hand side of the dual constraint is decreasing, first linearly with slope $2(z - y)$ until $t = \mu - y$, then convexly between $\mu - y$ and $\mu - z$, and finally linearly with slope $z - y$. The optimal dual solution is a straight line intersecting the constraint at the origin and at $t = \mu$. Accordingly, $\alpha_0 + \alpha_1\mu = (z - y)(2\mu - z/2 - y/2)/(2\mu)$, and the regret equals $(z - y)(p(2\mu - z/2 - y/2)/(2\mu) - c)$. The maximum regret, attained at $z = 2\mu(p - c)/p$, is then equal to $(py - 2\mu(p - c))^2/(4p\mu)$.

The robust order quantity minimizes the maximum of the following regrets:

$$\min_{y \geq 0} \max\{(py - 2\mu(p - c))^2/(4p\mu), (p - c)(\mu - y)\}$$

and is then equal to $y = 2\mu/p\sqrt{pc - c^2}$.

**Case 2:** $c/p \leq 1/2$.   When $z \geq y$, the right-hand side of the dual constraint is zero for $x \leq y$ and increasing after; it is convex between $y$ and $z$, and linearly increasing beyond $z$. The optimal dual solution is a straight line, intersecting the curve defined by the right-hand side at zero and $t = \mu$. Therefore, $\alpha_0 + \alpha_1\mu = (z - y)(2\mu - z/2 - y/2)/(2\mu)$, and the regret equals $(z - y)(p(2\mu - z/2 - y/2)/(2\mu) - c)$. The maximum regret, attained at $z = 2\mu(p - c)/p$, is then equal to $(py - 2\mu(p - c))^2/(4p\mu)$.

On the other hand, when $z \leq y$, the right-hand side is zero for $x \leq z$ and decreasing thereafter. The optimal dual solution is a horizontal line. Therefore, $\alpha_0 + \alpha_1\mu = 0$, and the regret equals $-c(z - y)$. The maximum regret, attained at $z = \mu$, is then equal to $-c(\mu - y)$.

The robust order quantity minimizes the following maximum regrets:

$$\min_{y \geq 0} \max\{(py - 2\mu(p - c))^2/(4p\mu), -c(\mu - y)\}$$

and is then equal to $y = 2\mu/p(p - \sqrt{pc - c^2})$. $\qquad\square$

# Proof of Theorem 7

When only the mean $\mu$ and the variance $\sigma^2$ are known, the dual problem (1.5) is the following:

$$\min \quad \alpha_0 + \alpha_1\mu + \alpha_2(\sigma^2 + \mu^2),$$

$$s.t. \quad \alpha_0 + \alpha_1 x + \alpha_2 x^2 \geq \min\{x, z\} - \min\{x, y\}, \quad \forall 0 \leq x.$$

Following Proposition 1, two cases must be considered: $z \geq y$ and $y \geq z$.

## Case 1: $z \geq y$.

The right hand sides of the constraints of the dual problem is piecewise linear increasing. A dual feasible function is any quadratic function $g(x)$ that, on the positive orthant, is nonnegative, lies above the line $x - y$ between $y$ and $z$, and above the line $z - y$ after $z$.

If the quadratic term of the function is zero, the problem reduces to finding a straight line, as in Theorem 1. Otherwise, in an optimal solution, either the parabola is tangent to $x - y$ at some point between $y$ and $z$ or it passes through the kink point $(z, y - z)$. Also, either its minimum value is zero on the half line $x \geq 0$, or it passes through the origin. Combining these possibilities gives rise to five different cases that we analyze next.

## Case 1.1: $g(x)$ is a straight line.
Following the proof of Theorem 1, two cases are possible. When $z \leq \mu$, the optimal solution is a horizontal line. The associated dual objective value is equal to $z - y$. On the other hand, when $z \geq \mu$, the optimal solution is a straight line passing through the origin. The optimal dual objective equals in this case $(z - y)\mu/z$.

## Case 1.2: $g(x)$ is tangent to $x - y$.
This analysis of this case is similar to the one performed by Bertsimas and Popescu (2002). If the quadratic function is tangent to $x - y$, it can be expressed as $g(x) = a(x - b)^2 + x - \alpha$ for some $a \geq 0$. The minimum

174

of the function, denoted by $x_0$, equals $b - 1/(2a)$.

If $b \geq 1/(2a)$, $x_0 \geq 0$. Since $g(x_0) = 0$, $a = 1/(b-y)$. After plugging this value of $a$ into the dual objective function, we minimize the function over all $b \in$ $[\max\{y, 1/(2a)\}, z]$. The minimum value equals $1/2(\mu - y) + (1/2)\sqrt{\sigma^2 + (\mu - y)^2}$ and is attained at $b = y + \sqrt{\sigma^2 + (\mu - y)^2}$. Trivially, $b \geq y$. On the other hand, $b \geq (1/2a)$ if and only if $y \geq (\mu^2 + \sigma^2)/(2\mu)$, and $b \leq z$ if and only if $y + \sqrt{\sigma^2 + (\mu - y)^2} \leq z$.

If $b \leq 1/(2a)$, $x_0 \leq 0$. Therefore, $g(x)$ must pass through the origin. Thus, $a = y/b^2$. After plugging this value of $a$ into the dual objective function, we minimize the function over all $b \in [y, \min\{1/(2a), z\}]$. The minimum value equals $\mu - y\mu^2/(\mu^2 + \sigma^2)$ and is attained at $b = (\sigma^2 + \mu^2)/\mu$. The condition that $b \leq (1/2a)$ simplifies to $y \leq (\mu^2 + \sigma^2)/(2\mu)$; therefore, the condition $b \leq y$ is automatically satisfied. Moreover, since $b \leq z$, $(\mu^2 + \sigma^2)/\mu \leq z$.

**Case 1.3: $g(x)$ passes through $(z, z - y)$.** Since $g(x)$ must lie above the piecewise linear function defined by the constraint, the derivative of $g(x)$ at the kink point must be less than 1, i.e., $2\alpha_2 z + \alpha_1 \leq 1$. The minimum of the function, denoted by $x_0$, equals $-\alpha_1/(2\alpha_2)$.

If $-\alpha_1/(2\alpha_2) \leq 0$, the quadratic function passes through the origin, i.e. $g(0) = 0$. Accordingly, the dual objective function simplifies to $\alpha_2(\sigma^2 + \mu^2 - \mu z) + \mu(z - y)/z$. Minimizing the objective with respect to $\alpha_2$ gives rise to three possible cases. When $z \leq (\sigma^2 + \mu^2)/\mu$, the minimum is attained at $\alpha_2 = 0$ (straight line) and equals $\mu(z - y)/z$. If $z \geq (\sigma^2 + \mu^2)/\mu$ and $z \leq 2y$, the minimum equals $(\sigma^2 + \mu^2)(z - y)/z^2$, and $x_0 = 0$. If $z \geq (\sigma^2 + \mu^2)/\mu$ and $z \geq 2y$, the minimum equals $(\sigma^2 + \mu^2)y/z^2 + (z - 2y)/z$ with $g'(z) = 1$. Clearly, these last two cases are "border situations", resulting from the distinction between cases, and will never be optimal.

On the other hand, if $-\alpha_1/(2\alpha_2) \geq 0$, the quadratic function is minimized at some $x_0 \geq 0$. Thus, $g(x)$ is assumed to pass through $(z, z - y)$ with a derivative less than 1, and to have a minimum value of zero, attained on the interval $[0, y]$. From these conditions, the dual feasible function can be expressed as a function of $\alpha_2$ only. Consequently, the dual objective function can also be expressed as the following

function of $\alpha_2$: $\alpha_2(\sigma^2 + \mu^2 - 2\mu z + z^2) + 2\sqrt{\alpha_2}\sqrt{z-y}(\mu - z) + z - y$. Minimizing the dual objective over all nonnegative values of $\alpha_2$, such that the above conditions are met, gives rise to the following cases. If $z \leq (\sigma^2 + \mu^2)/\mu$, the minimum dual objective value equals $(\sigma^2 + \mu^2)(z-y)/z^2$, with $x_0 = 0$. If $z^2/2 + \mu y - zy \geq (\sigma^2 + \mu^2)/2$, the dual objective value equals $\mu - y + (\sigma^2 + (\mu - z)^2)/(4(z-y))$, with $g'(z = 1$. Clearly, these two cases are "border situations", resulting from the distinction between cases, and will never be optimal. Finally, if $z \geq (\sigma^2 + \mu^2)/\mu$ and if $z^2/2 + \mu y - zy \leq (\sigma^2 + \mu^2)/2$, the minimum dual objective value equals $\sigma^2(z - y)/(\sigma^2 + (\mu - z)^2)$, attained when $\sqrt{\alpha_2} = (z - \mu)\sqrt{z-y}/(\sigma^2 + (\mu - z)^2)$.

**Case 1 Summary.** All border situations can be discarded from consideration, as they arose from the separation between cases. Therefore, the optimal value of the dual problem (1.5) is equal to:

$$
\begin{array}{ll}
z - y, & \text{if } \mu \geq z, \\[4pt]
\frac{\mu}{z}(z - y) & \text{if } \mu \leq z \leq \frac{\sigma^2 + \mu^2}{\mu}, \\[4pt]
(z - y)\frac{\sigma^2}{\sigma^2 + (z-\mu)^2} & \text{if } \frac{\sigma^2 + \mu^2}{\mu} \leq z \leq y + \sqrt{\sigma^2 + (y - \mu)^2}, \\[4pt]
\mu - y\frac{\mu^2}{\sigma^2 + \mu^2} & \text{if } \frac{\sigma^2 + \mu^2}{\mu} \leq z, \text{ and } y \leq \frac{\sigma^2 + \mu^2}{2\mu}, \\[4pt]
(1/2)(\mu - y) + (1/2)\sqrt{\sigma^2 + (\mu - y)^2} & \text{if } z \geq y + \sqrt{\sigma^2 + (y - \mu)^2}, \text{ and } y \leq \frac{\sigma^2 + \mu^2}{2\mu}.
\end{array}
$$

The two candidate solutions obtained in Case 1.2 can be disregarded. In fact, they are independent of $z$. Therefore, the regret associated with these solutions equals a constant minus $c(z - y)$. Maximizing the regret over all values of $z$ will take $z$ as small as possible. Therefore, the regret associated with the last two solutions is dominated by $\sigma^2(z - y)/(\sigma^2 + (\mu - z)^2)$.

Similarly, the solution $z - y$ can also be discarded from consideration. The associated regret, $(p - c)(z - y)$ is increasing in $z$. Taking the largest $z$ makes $z - y$ equal to $\mu/z(z - y)$.

As a result, the maximum regret must be maximized between

$$(p\frac{\mu}{z} - c)(z - y) \qquad \text{if } \mu \leq z \leq \frac{\sigma^2 + \mu^2}{\mu}, \text{and}$$

$$(p\frac{\sigma^2}{\sigma^2 + (z-\mu)^2} - c)(z - y) \quad \text{if } \frac{\sigma^2 + \mu^2}{\mu} \leq z \leq y + \sqrt{\sigma^2 + (y - \mu)^2}.$$

The first term is a concave function of $z$, while the second term is concave-convex with only one maximum on its interval.

## Case 2: $z \leq y$.

The right hand sides of the constraints of the dual problem is piecewise linear decreasing. A dual feasible function is any quadratic function $g(x)$ that, on the positive orthant, is nonnegative, lies above the line $x - z$ between $z$ and $y$, and above the line $z - y$ after $y$.

If the quadratic term of the function is zero, the problem reduces to finding a straight line, as in Theorem 1. Otherwise, in an optimal solution, either the parabola is tangent to $x - z$ at some point between $z$ and $y$ or it passes through the kink point $(z, 0)$. In the last two cases, its minimum value is $z - y$ on the half line $x \geq y$. Combining these possibilities gives rise to three different cases that we analyze next.

**Case 2.1: $g(x)$ is a straight line.** Since the constraint right-hand side are decreasing, a horizontal line at zero is a candidate solution. The optimal dual objective equals zero in this case.

**Case 2.2: $g(x)$ is tangent to $x - z$.** In this case, the function can be expressed as $g(x) = a(x - b)^2 + z - x$, for some $a \geq 0$. The minimum of the function, denoted by $x_0$, is then equal to $b + 1/(2a)$. Since $g(x_0) = z - y$, $b = y - 1/(4a)$. The dual objective function can then be expressed as a function of $a$ only, namely $a(\sigma^2 + \mu^2) + 2a(1/(4a) - y)\mu + z - \mu + a(1/(4a) - y)^2$. When minimized over all nonnegative $a$, the function attains its minimum at $a = 1/(4\sqrt{\sigma^2 + (\mu - y)^2})$ and is equal to $(1/2)\sqrt{\sigma^2 + (\mu - y)^2} + (\mu - y)/2 + z - \mu$. The point of tangency between the quadratic function and the constraint is equal to $b = y - \sqrt{\sigma^2 + (\mu - y)^2}$ is trivially less than

177

$y$, as well as the minimum point $x_0$. It is greater than $z$ if $z \leq y - \sqrt{\sigma^2 + (\mu - y)^2}$.

**Case 2.3:** $g(x)$ **passes through** $(z, 0)$. Since $g(x)$ must lie above the piecewise linear function defined by the constraint, the derivative of $g(x)$ at the kink point must be greater than -1, i.e., $2\alpha_2 z + \alpha_1 \geq -1$. The minimum of the function, denoted by $x_0$, equals $-\alpha_1/(2\alpha_2)$. Thus, $g(x)$ is assumed to pass through $(z, 0)$ with a derivative greater than -1, and to have a minimum value of $z - y$, attained on the interval $[y, \infty)$. From these conditions, the dual feasible function can be expressed as a function of $\alpha_2$ only. Consequently, the dual objective function can also be expressed as the following function of $\alpha_2$: $\alpha_2(\sigma^2 + \mu^2 - 2\mu z + z^2) + 2\sqrt{\alpha_2}\sqrt{z - y}(\mu - z)$. Minimizing the dual objective over all nonnegative values of $\alpha_2$, such that the above conditions are met, gives rise to the following cases. If $z \geq \mu$, the derivative of the objective function with respect to $\alpha_2$ is always positive. The optimal solution is to take $\alpha_2$ as small as possible, i.e., $\alpha_2 = 0$. In this case, the dual optimal solution is a horizontal line at zero, similarly to Case 2.1. If $z \leq \mu$ and $2(\mu - z)(z - y) \leq \sigma^2 + (\mu - z)^2$, the minimum objective is attained when $\sqrt{\alpha_2} = (\mu - z)\sqrt{y - z}/(\sigma^2 + (\mu - z)^2)$, and equals $(z - y)(\mu - z)^2/(\sigma^2 + (\mu - z)^2)$. The other cases are "border situations", arising from the distinction between cases in the analysis, and can be discarded from future consideration.

**Case 2 Summary.** All border situations can be discarded from consideration, as they arose from the separation between cases. Therefore, the optimal value of the dual problem (1.5) is equal to:

$$
\begin{cases}
0, & \text{if } \mu \geq z, \\
(z - y)\frac{(\mu - z)^2}{\sigma^2 + (z - \mu)^2} & \text{if } \mu \geq z \geq y - \sqrt{\sigma^2 + (y - \mu)^2}, \\
(1/2)(\mu - y) + (1/2)\sqrt{\sigma^2 + (\mu - y)^2} + z - \mu & \text{if } z \leq y - \sqrt{\sigma^2 + (y - \mu)^2}.
\end{cases}
$$

Optimizing the regret over $z$, it turns out that only the second solution is relevant, as the derivative with respect to $z$ of the first regret is always negative and that of the third regret is always positive on its interval. Furthermore, one can show that the

178

second regret function is concave on its interval of definition. Therefore, the optimal regret is

$$\max_{\min\{\mu,y\}\geq z\geq y-\sqrt{\sigma^2+(\mu-y)^2}}\{(z-y)(p\frac{(z-\mu)^2}{\sigma^2+(z-\mu)^2}-c)\}.$$

Equating both regrets gives rise to the theorem statement. $\qquad\square$

## Justification of Approximation 1

First, we bound from above the left-hand side of (1.8) with $\max_x(p\frac{\sigma^2}{\sigma^2+(x-\mu)^2}-c)(x-y)$ in which $\max\{\mu,y\}\leq x\leq y+\sqrt{\sigma^2+(y-\mu)^2}$. This is indeed an upper bound, since the function $(p\frac{\mu}{x}-c)(x-y)$ is the optimal value of the dual problem (1.5)–a minimization problem–when $x\leq(\sigma^2+\mu^2)/\mu$.

From the proof of Theorem 7, the maximum of the function $(x-y)(p\sigma^2/(\sigma^2+(\mu-x)^2)-c)$ is the maximum regret from not ordering $x$, when $x\geq y$. In fact, the worst-case demand distribution is a two-point distribution, with a probability mass of $\sigma^2/(\sigma^2+(\mu-x)^2)$ at $x$, corresponding to the right-hand side of Chebyshev's inequality (see, e.g., Bertsimas and Popescu 2005). By optimality of $x$, we have that $\sigma^2/(\sigma^2+(\mu-x)^2)\geq c/p$. Solving the inequality for $x$ gives rise to the following upper bound on $x$: $\mu+\sigma\sqrt{(p-c)/p}$. From (1.8), $\mu$ is a lower bound of $x$, when $x\geq y$. Therefore, rather than maximizing the regret of ordering too little over all possible values of $x$, we evaluate the function at the middle point $x=\mu+(1/2)\sigma\sqrt{(p-c)/c}$.

Similarly, the maximum of the function $(x-y)(p\frac{(x-\mu)^2}{\sigma^2+(x-\mu)^2}-c)$ is the maximum regret from ordering too much, when $x\leq y$. The worst-case demand distribution is also a two-point distribution, with a probability mass of $\sigma^2/(\sigma^2+(\mu-x)^2)$ at $x$. By optimality of $x$, we must have that $(x-\mu)^2/(\sigma^2+(\mu-x)^2)\leq c/p$. Solving the inequality gives rise to the following lower bound of $x$: $\mu-\sigma\sqrt{c/(p-c)}$. From (1.8), $\mu$ is an upper bound of $x$, when $x\leq y$. Therefore, rather than maximizing the regret of ordering too much over all possible values of $x$, we evaluate the function at the middle point $x=\mu-(1/2)\sigma\sqrt{c/(p-c)}$.

Equating both sides, and solving for $y$, we obtain (1.9).

By construction, the approximation is not expected to be accurate when the maximizing term in the left-hand side of (1.8) is less than $(\sigma^2 + \mu^2)/\mu$, since the associated function $(p\frac{\sigma^2}{\sigma^2 + (x-\mu)^2} - c)(x - y)$ is replaced with $(x - y)(px/\mu - c)$. In particular, this will happen when the upper bound on $x$, $\mu + \sigma\sqrt{(p-c)/p}$, is less than $(\sigma^2 + \mu^2)/\mu$, i.e., when $\sigma/\mu \geq \sqrt{(p-c)/p}$. Accordingly, the approximation is not expected to be accurate for high coefficients of variation $\sigma/\mu$ or low profit margins $p/c$. $\quad\square$

## Proof of Theorem 8

The minimization problem in (1.14) consists in minimizing a concave function (minimum of linear functions) over a polyhedron. Therefore, an optimal solution is an extreme point of the polyhedron. Since $\Gamma$ is integer, every extreme point of the polyhedron $\{\boldsymbol{\delta} : 0 \leq \boldsymbol{\delta} \leq 1, \mathbf{1}'\boldsymbol{\delta} \leq \Gamma\}$ different from zero has $\Gamma$ coordinates equal to 1 and $M - \Gamma$ coordinates equal to zero. Since $\delta_j$ can only take two values in an optimal solution, the inner problem is equivalent to:

$$
\begin{aligned}
\min \quad & \mathbf{p}' \min\{\mathbf{M}, \mathbf{Y}\}(1 - \boldsymbol{\delta}) + \mathbf{p}' \min\{\mathbf{M} - \mathbf{S}, \mathbf{Y}\}\boldsymbol{\delta}, \\
s.t. \quad & \mathbf{1}'\boldsymbol{\delta} \leq \Gamma, \\
& 0 \leq \boldsymbol{\delta} \leq 1,
\end{aligned}
$$

where $\mathbf{M}, \mathbf{S}$, and $\mathbf{Y}$ are diagonal matrices whose diagonals are equal to $\mathbf{m}, \mathbf{s}$, and $\mathbf{y}$ respectively.

Associating a dual variable $\alpha$ with the budget constraint and vector of dual variables $\boldsymbol{\beta}$ with the upper bound constraint, the dual problem can be formulated as follows:

$$
\begin{aligned}
\max \quad & \alpha\Gamma + \boldsymbol{\beta}'\mathbf{1} + \mathbf{p}' \min\{\mathbf{y}, \mathbf{m}\} \\
s.t. \quad & \alpha\mathbf{1} + \boldsymbol{\beta} \leq (\min\{\mathbf{M} - \mathbf{S}, \mathbf{Y}\} - \min\{\mathbf{M}, \mathbf{Y}\})\mathbf{p}, \\
& \alpha \leq 0, \boldsymbol{\beta} \leq 0.
\end{aligned}
$$

The dual problem is a maximization problem. Therefore, the maximin network

RM problem (1.14) simplifies to the following maximization problem:

$$\max \qquad \alpha\Gamma + \beta'\mathbf{1} + \mathbf{p}'\min\{\mathbf{y}, \mathbf{m}\}$$

$$s.t. \qquad \mathbf{Ay} \leq \mathbf{b},$$

$$\alpha\mathbf{1} + \beta + \min\{\mathbf{M}, \mathbf{Y}\}\mathbf{p} \leq \mathbf{Yp},$$

$$\alpha\mathbf{1} + \beta + \min\{\mathbf{M}, \mathbf{Y}\}\mathbf{p} \leq (\mathbf{M} - \mathbf{S})\mathbf{p},$$

$$\alpha \leq 0, \beta \leq 0, \mathbf{y} \geq \mathbf{0}.$$

Introducing the variable $w_j = \min\{y_j, m_j\}$ simplifies the problem to a linear optimization problem. $\qquad\qquad\square$

## Proof of Lemma 1

(a) The minimax regret problem can be formulated as a moment bound problem, with only one constraint per distribution, namely $\int_{l_j}^{u_j} P(x)dx = 1$. By linear programming theory, there exists an optimal solution which is an extreme point of the polyhedron. Since there is only one constraint, an extreme point has only one positive variable $P(x)$, corresponding to a unit impulse distribution.

(b) If $z_j \geq l_j$ (resp. $z_j < l_j$), it is never optimal to have a demand $d_j > z_j$ (resp. $d_j > l_j$). Otherwise, $d_j$ could be reduced without diminishing the revenue associated with $z_j$ but with potentially reducing the revenue associated with $y_j$, leading to a regret decrease.

(c) First, suppose that $y_j \leq l_j$. If $z_j \geq l_j$, then $d_j = z_j$ by (b), and the regret equals to $p_j(z_j - y_j)$; if $z_j < l_j$, $d_j = l_j$ and the regret might be positive or negative and equals $p_j(z_j - y_j)$. Second, suppose that $y_j \geq l_j$. If $z_j \leq l_j$, then $d_j = l_j$ and the regret is equal to $p_j(z_j - l_j)$. On the other hand, if $z_j \geq l_j$ and $z_j \leq y_j$, $d_j = z_j$ and the regret is zero. Finally, if $z_j \geq y_j \geq l_j$, $d_j = z_j$ and the regret equals $p_j(z_j - y_j)$. $\qquad\qquad\square$

# Proof of Theorem 9

Let $\delta_j^0, \delta_j^l$, and $\delta_j^u$ the probabilities that the optimal booking limit $Z_j$ equals 0, $l_j$, and $u_j$ respectively. Since any feasible value for $Z_j$ can be expressed as a convex combination of these three points, $\delta_j^0 + \delta_j^l + \delta_j^u = 1$, and these probabilities are between 0 and 1. From Lemma 1, the regret equals $p_j(u_j - y_j)$ when $z_j = u_j$, $p_j \max\{0, l_j - y_j\}$ when $z_j = l_j$, and $-p_j \min\{y_j, l_j\}$ when $z_j = 0$. Under Assumption 1, the objective function can also be expressed as a convex combination of these three points. Accordingly, the regret maximization problem can be rewritten as

$$\max \quad \mathbf{p}'(\mathbf{U} - \mathbf{Y})\delta^u + \mathbf{p}' \max\{0, \mathbf{L} - \mathbf{Y}\}\delta^l - \mathbf{p}' \min\{\mathbf{L}, \mathbf{Y}\}\delta^0,$$

$$s.t. \qquad \delta^u + \delta^l + \delta^0 = 1,$$

$$\mathbf{A}(\mathbf{U}\delta^u + \mathbf{L}\delta^l) \leq \mathbf{b},$$

$$\delta^u, \delta^l, \delta^0 \geq 0.$$

We associate a vector of dual variables $\boldsymbol{\pi}$ with the capacity constraints and a vector of dual variables $\mathbf{q}$ with the constraints that probabilities sum up to unity. By strong duality, the above problem is equivalent to the following dual problem:

$$\min \qquad \boldsymbol{\pi}'\mathbf{b} + \mathbf{q}'\mathbf{1}$$

$$s.t. \quad \boldsymbol{\pi}'\mathbf{AU} + \mathbf{q}' \geq \mathbf{p}'(\mathbf{U} - \mathbf{Y}),$$

$$\boldsymbol{\pi}'\mathbf{AL} + \mathbf{q}' \geq \mathbf{p}' \max\{0, \mathbf{L} - \mathbf{Y}\},$$

$$\mathbf{q}' \geq -\mathbf{p}' \max\{\mathbf{Y}, \mathbf{L}\},$$

$$\boldsymbol{\pi} \geq 0.$$

Plugging this inner problem into the general minimax regret problem, we obtain the theorem statement. $\qquad\qquad\square$

## A.2  Chapter 2

### Lemma 4

**Lemma 4.** *Fix $\bar{F}(y) = \varphi$. Then,*

$$\inf_{F \in \mathcal{D}} \int_0^y \bar{F}(x) = y\varphi,$$

*where $\mathcal{D}$ is the set of all IGFR distributions; the minimizing distribution is the following:*

$$\bar{F}(x) = \begin{cases} 1 & \text{if } 0 \leq x \leq y\varphi^{1/l}, \\ \varphi y^l x^{-l} & \text{if } y\varphi^{1/l} \leq x \leq y. \end{cases}$$

*when $l$ tends to zero.*

*Proof.* With the proposed distribution,

$$\begin{aligned} \int_0^y \bar{F}(x)dx &= y\varphi^{1/l} + \varphi y/(1-l) - y\varphi^{1/l}/(1-l), \\ &= \varphi y, \end{aligned}$$

when $l$ tends to zero. The distribution has a constant generalized failure rate, equal to $l$, between $y\varphi^{1/l}$ and $y$, and is thus IGFR. $\qquad \square$

### Lemma 5

**Lemma 5.** *Fix $\bar{F}(y) = \varphi$ and $g(y) = k$. Then, for $x \geq y$,*

$$\max_{F \in \mathcal{D}} \bar{F}(x) = x^{-k}y^k\varphi,$$

*where $\mathcal{D}$ is the set of all IGFR distributions.*

*Proof.* Expressing $\bar{F}(x)$ as an exponential function of its hazard rate (Barlow and

Proschan 1965), we have the following

$$
\begin{aligned}
\bar{F}(x) &= \exp(-\int_0^x h(\xi)d\xi), \\
&= \bar{F}(y)\exp(-\int_y^x h(\xi)d\xi), \\
&\leq \bar{F}(y)\exp(-\int_y^x (k/\xi)d\xi), \\
&= \varphi y^k x^{-k},
\end{aligned}
$$

where the inequality comes from the IGFR property of $F(x)$. $\qquad\square$

## Lemma 6

**Lemma 6.** *Fix $\bar{F}(y^d) = \varphi$ and $g(y^d) = k$. Then*

$$
\sup_{D \in \mathcal{D}} \frac{-cy^c + pE[\min\{y^c, D\}]}{-cy^d + pE[\min\{y^d, D\}]} = \frac{k}{1-k}\frac{r(\frac{\varphi}{r})^{1/k} - \varphi}{\varphi - r}, \tag{A.2}
$$

*where $\mathcal{D}$ is the set of IGFR distributions. The maximizing distribution is piecewise Pareto with a breakpoint at $y^d$:*

$$
\bar{F}(x) = \begin{cases} x^{-l}(y^d)^l \varphi, & \text{for } y^d(\varphi)^{1/l} \leq x < y^d, \\ x^{-k}(y^d)^k \varphi, & \text{for } x \geq y^d, \end{cases}
$$

*when $l$ tends to zero.*

*Proof.* From integrating by parts, $E[\min\{y, D\}] = \int_0^y \bar{F}(x)dx$. Hence, after dividing both the numerator and the denominator by $p$, the ratio of profits can be written as follows:

$$
\frac{-cy^c + pE[\min\{y^c, D\}]}{-cy^d + pE[\min\{y^d, D\}]} = \frac{-ry^c + \int_0^{y^d} \bar{F}(x)dx + \int_{y^d}^{y^c} \bar{F}(x)dx}{-ry^d + \int_0^{y^d} \bar{F}(x)dx}.
$$

Since the ratio is greater than 1, it is maximized when $\int_0^{y^d} \bar{F}(x)dx$ is the smallest. From Lemma 4, $\int_0^{y^d} \bar{F}(x)dx \geq y^d \bar{F}(y^d) = y^d \varphi$. Replacing $\int_0^{y^d} \bar{F}(x)dx$ with its lower

184

bound gives rise to the following bound:

$$\frac{-ry^c + y^d\varphi + \int_{y^d}^{y^c} \bar{F}(x)dx}{-ry^d + \varphi y^d}$$

The ratio is maximized when $\int_{y^d}^{y^c} \bar{F}(x)dx$ is the largest. Using Lemma 5 gives rise to the following inequality:

$$\int_{y^d}^{y^c} \bar{F}(x)dx \leq (y^d)^k \varphi \int_{y^d}^{y^c} x^{-k}dx,$$

$$= \frac{\varphi}{1-k}(y^c(\frac{y^d}{y^c})^k - y^d).$$

After replacing $\int_{y^d}^{y^c} \bar{F}(x)dx$ by its upper bound, and dividing both sides of the ratio by $y^d$, one obtains the following bound:

$$\frac{-r\frac{y^c}{y^d} + \varphi + \varphi\frac{1}{1-k}((\frac{y^c}{y^d})^{1-k} - 1)}{-r + \varphi}.$$

The ratio is increasing with $y^c/y^d$. Using Lemma 5, $r = \bar{F}(y^c) \leq (\frac{y^d}{y^c})^k \varphi$. Thus, the ratio is maximized when $y^c/y^d = (\varphi/r)^{1/k}$, giving rise to the lemma statement. $\square$

## Proof of Theorem 10

The proof is based on Lemmas 4-6 in the appendix. Suppose $y^d$ is an interior solution and satisfies (2.2), i.e., $\bar{F}(y^d) = r/(1-k)$, with $k = g(y^d)$. Replacing $\varphi$ with $r/(1-k)$ in (A.2) simplifies the Price of Anarchy to:

$$\frac{-cy^c + pE[\min\{y^c, D\}]}{-cy^d + pE[\min\{y^d, D\}]} \leq (1-k)^{-1} - (1-k)^{-1/k}.$$

The right-hand side is an increasing function of $k$ and attains its maximum when $k = 1 - r$ (since $r/(1-k) = \bar{F}(y^d) \leq 1$), leading to the theorem statement. Notice that, since $\bar{F}(y^d) = 1$ in the worst-case, $y^d$ is not in the interior but corresponds to the lowest value of the support of the distribution. $\square$

185

## On Remark 1

If $\pi \geq r$, the optimal inventory level in the integrated supply chain, $y^c$, equals $u$. In the decentralized supply chain, the manufacturer can induce the retailer to order $l$ or $u$. If she induces the retailer to order $l$, she will offer a wholesale price $p$ and obtain a profit equal to $(p - c)l$. On the other hand, if she induces the retailer to order $u$, the maximum wholesale price she can offer is $\pi p$, which would leave her with a profit equal to $(\pi p - c)u$. If $\pi \leq r + (1 - r)l/u$, she will choose a price $p$ and induce an order quantity $l$. The ratio of expected profits for the entire supply chain is then equal to:

$$\frac{E[\Pi(y^c, D)]}{E[\Pi(y^d, D)]} = \frac{p(\pi u + (1 - \pi)l) - cu}{(p - c)l}.$$

The ratio increases with $\pi$ and is maximized when $\pi = r + (1 - r)l/u$. At its maximum, the ratio equals $2 - l/u$. □

## Proof of Lemma 2

The proof is by induction, by moving upstream along the supply chain. With only two stages, the lemma reduces to (2.2).

For the induction, assume that if the supply chain consists of $i$ stages and has an inbound cost $w$, the order quantity satisfies $\bar{F}(y)(1 - g(y))^{i-1} \geq w/p$. Consider now the stage upstream of that supply chain, i.e., stage $i + 1$. With an inbound cost of materials $c$, stage $i + 1$ will choose her transfer price $w$ so as to maximize her profits, anticipating the order quantity of the downstream stages, that is,

$$\max_{w,y} \quad (w - c)y,$$
$$s.t. \quad \bar{F}(y)(1 - g(y))^{i-1} \geq w/p.$$

Since the stage wants to maximize its profits, it will offer $w = p\bar{F}(y)(1 - l(y))^{i-1}$, and choose the induced level of inventory $y$ to maximize $(p\bar{F}(y)(1 - g(y))^{i-1} - c)y$. Differentiating the profit function with respect to $y$, and using the fact that $g(y)$ is increasing (because of the IGFR property), we obtain that $p\bar{F}(y)(1 - g(y))^{i} - c \geq 0$,

completing the induction step. □

## Proof of Theorem 11

The proof is based on Lemmas 4-6 in the appendix. The bound (A.2) in Lemma 6 is decreasing with $\varphi$. From Lemma 3, $\varphi \geq r/(1-g(y))^{n-1}$. Replacing $\varphi$ with $r/(1-k)^{n-1}$ in (A.2), we obtain an increasing function of $k$, attaining its maximum when $\bar{F}(y^d) = 1$, i.e., when $k = 1 - r^{1/(n-1)}$. Notice that since the worst-case distribution is Pareto, the inequality in Lemma 2 is tight in the worst case. □

## Proof of Theorem 12

The proof is based on Lemmas 4-6 in the appendix. In Lemma 6, the worst-case distribution is such that $\int_0^{y^d} \bar{F}(x)dx = \bar{F}(y^d)y^d$. Thus, optimality condition (2.3) simplifies to $\bar{F}(y^d) = r(1+k)$, where $k = g(y^d)$. Replacing $\varphi$ with $r(1+k)$ into upper bound (A.2) simplifies the Price of Anarchy to $((1+k)^{1/k} - 1 - k)/(1-k)$, which is decreasing with $k$ and is maximized when $k$ tends to zero. □

## Proof of Lemma 3

The proof is by induction, by moving downstream along the supply chain. With only two stages, the claim holds at equality by (2.3).

For the induction, assume that, if the supply stage consists of $i$ stages and has a unit selling price $w$, the level of inventory in the supply chain satisfies $\bar{F}(y) \geq (c/w)(1 + \ell(y))^{i-1}$. Consider now the stage downstream of the supply chain, i.e., stage $i + 1$. With a unit selling price $p$, stage $i + 1$ will choose a transfer price $w$ in order to maximize her profits, anticipating the level of inventory in the upstream supply chain, that is,

$$\max_{w,y} \quad (p - w) \int_0^y \bar{F}(x)dx,$$
$$s.t. \quad \bar{F}(y) \geq \frac{c}{w}(1 + \ell(y))^{i-1}.$$

187

Since the stage wants to maximize its profits, it will offer $w = c/\bar{F}(y)(1 + \ell(y))^{i-1}$, and choose the induced level of inventory $y$ to maximize $(p - c/\bar{F}(y)(1 + \ell(y))^{i-1}) \int_0^y \bar{F}(x)dx$. Differentiating the profit function with respect to $y$, and using the fact that $\ell(y)$ is increasing (because of the IGFR property), we obtain that $p\bar{F}(y) - c(1 + \ell(y))^i \geq 0$, completing the induction step. $\qquad \square$

## Proof of Theorem 13

The proof is based on Lemmas 4-6 in the appendix. The bound (A.2) in Lemma 6 is decreasing with $\varphi$. From Lemma 3, $\varphi \geq r(1 + \ell(y))^{n-1}$. With the distribution that attains the upper bound (A.2), $\int_0^{y^d} \bar{F}(x)dx \rightarrow y^d\bar{F}(y^d)$; accordingly, $\ell(y^d) \rightarrow g(y^d)$. Replacing $\varphi$ with $r(1 + k)^{n-1}$ in (A.2), we obtain a decreasing function of $k$, which is maximized when $k$ tends to 0. Since $k \rightarrow 0$, the inequality in Lemma 3 is tight with the worst-case distribution. $\qquad \square$

## Proof of Theorem 14

The proof is similar to the proof of Theorem 10. Using (2.4), we replace $\varphi$ with $r/(1 - kn)$ in (A.2) of Lemma 6. The resulting upper bound is increasing with $k$ and maximized when $k = (1 - r)/n$ (since $\bar{F}(y^d) = r/(1 - kn) \leq 1$), leading to the theorem statement. $\qquad \square$

## Lemma 7

**Lemma 7.** *If the demand distribution is IGFR, is twice differentiable, and is such that $f'(x)\bar{F}(x)/(f(x))^2$ is increasing, the function*

$$\sum_{i=1}^n (p - \frac{c}{\bar{F}(\bar{y}_i)}) \int_{\bar{y}_{i-1}}^{\bar{y}_i} \bar{F}(x)dx,$$

*where $\bar{y}_0 = 0$, is concave in $\bar{y}_1, ..., \bar{y}_n$.*

*Proof.* Let us decompose the function as follows

$$\sum_{i=1}^{n}\left(p - \frac{c}{\bar{F}(\bar{y}_i)}\right)\int_{\bar{y}_{i-1}}^{\bar{y}_i}\bar{F}(x)dx = \sum_{i=1}^{n}p\int_{0}^{\bar{y}_n}\bar{F}(x)dx - \sum_{i=1}^{n}\frac{c}{\bar{F}(\bar{y}_i)}\int_{\bar{y}_{i-1}}^{\bar{y}_i}\bar{F}(x)dx.$$

The first term is a concave function of $\bar{y}_n$. A sufficient condition for the second term to be concave is that each term $\int_{\bar{y}_{i-1}}^{\bar{y}_i}\bar{F}(x)dx/\bar{F}(\bar{y}_i)$ is convex in $\bar{y}_{i-1}$ and $\bar{y}_i$, as the sum of convex functions is convex and the negative of a convex function is concave. Since the distribution function is twice differentiable, the function is convex if and only if the Hessian is positive semidefinite. Taking the second derivatives gives the following Hessian:

$$\begin{bmatrix} \frac{f(\bar{y}_{i-1})}{\bar{F}(\bar{y}_i)} & -\frac{f(\bar{y}_i)\bar{F}(\bar{y}_{i-1})}{(\bar{F}(\bar{y}_i))^2} \\ -\frac{f(\bar{y}_i)\bar{F}(\bar{y}_{i-1})}{(\bar{F}(\bar{y}_i))^2} & \frac{f(\bar{y}_i)}{\bar{F}(\bar{y}_i)} + \int_{\bar{y}_{i-1}}^{\bar{y}_i}\bar{F}(x)dx\left(\frac{f'(\bar{y}_i)}{(\bar{F}(\bar{y}_i))^2} + 2\frac{(f(\bar{y}_i))^2}{(\bar{F}(\bar{y}_i))^3}\right) \end{bmatrix}.$$

Using a similar proof as Lemma 1 in Cachon (2004), one can show that the quantity $f(x)/(\bar{F}(x))^2\int_{\alpha}^{x}\bar{F}(\xi)d\xi$ is increasing in $x$ if the demand distribution is IGFR. (The original lemma was shown for the case $\alpha = 0$.) Accordingly, the diagonal elements of the Hessian are nonnegative.

To complete the proof, we need to show that the determinant of the Hessian is nonnegative. It is equal to zero when $\bar{y}_{i-1} = \bar{y}_i$. Suppose now that it is equal to zero for some $\bar{y}_{i-1}$ and let us show that the derivative with respect to $\bar{y}_{i-1}$ is negative. If this is true, then the determinant is nonnegative.

The derivative of the determinant with respect to $\bar{y}_{i-1}$ is equal to

$$\int_{\bar{y}_{i-1}}^{\bar{y}_i}\bar{F}(x)dx\left(\frac{f'(\bar{y}_i)f'(\bar{y}_{i-1})}{(\bar{F}(\bar{y}_i))^3} + 2\frac{(f(\bar{y}_i))^2 f'(\bar{y}_{i-1})}{(\bar{F}(\bar{y}_i))^4}\right)$$
$$- \frac{f'(\bar{y}_i)f(\bar{y}_{i-1})\bar{F}(\bar{y}_{i-1})}{(\bar{F}(\bar{y}_i))^3} + \frac{f(\bar{y}_i)f'(\bar{y}_{i-1})}{(\bar{F}(\bar{y}_i))^2}.$$

If the determinant is equal to zero, then by rearranging terms,

$$\int_{\bar{y}_{i-1}}^{\bar{y}_i}\bar{F}(x)dx = \frac{-f(\bar{y}_i)f(\bar{y}_{i-1})(\bar{F}(\bar{y}_i))^2 + (\bar{F}(\bar{y}_{i-1}))^2(f(\bar{y}_i))^2}{f'(\bar{y}_i)f(\bar{y}_{i-1})\bar{F}(\bar{y}_i) + 2f(\bar{y}_{i-1})(f(\bar{y}_i))^2}.$$

189

Replacing $\int_{\bar{y}_{i-1}}^{\bar{y}_i} \bar{F}(x)dx$ with the above function simplifies the derivative of the determinant to

$$\frac{(\bar{F}(\bar{y}_{i-1}))^2}{(\bar{F}(\bar{y}_i))^3} f(\bar{y}_i) \left( \frac{f'(\bar{y}_{i-1})f(\bar{y}_i)}{f(\bar{y}_{i-1})\bar{F}(\bar{y}_i)} - \frac{f'(\bar{y}_i)f(\bar{y}_{i-1})}{f(\bar{y}_i)\bar{F}(\bar{y}_{i-1})} \right),$$

which is nonpositive if $f'(x)\bar{F}(x)/(f(x))^2$ is nondecreasing, since $\bar{y}_{i-1} \leq \bar{y}_i$. $\qquad \square$

## Proof of Theorem 15

Lemma 4 showed that the distribution that minimizes $\int_0^{\bar{y}_n^d} \bar{F}(x)dx$ for fixed $\bar{y}_n^d$ and $\bar{F}(\bar{y}_n^d)$, is Pareto with a generalized failure rate converging to zero. In this proof, we generalize this result by considering a piecewise Pareto distribution, with breakpoints at $\bar{y}_1^d, ..., \bar{y}_n^d$. The generalized failure rates of all pieces, except the last one, tend to zero, and the convergence to zero is faster for the pieces associated with lower values. Since this distribution generalizes the result obtained in Lemma 4, it is clear that $\int_0^{\bar{y}_i^d} \bar{F}(x)dx = \bar{y}_i^d \bar{F}(\bar{y}_i^d)$, for $i = 1, ..., n$.

In addition, we will show by induction that with that distribution, $\bar{F}(\bar{y}_{i-1}^d) = \bar{F}(\bar{y}_i^d)(1 + g(\bar{y}_{i-1}^d)\mathcal{E}^{i-2}(1))$ for $i = 2, ..., n$. With that result, $\bar{F}(\bar{y}_n^d) = r(1 + g(\bar{y}_n^d)\mathcal{E}^{n-1}(1))$. Replacing $\varphi$ with $r(1 + k\mathcal{E}^{n-1}(1))$ in upper bound (A.2) in Lemma 6 simplifies the Price of Anarchy to a function of $k$ only. The function is decreasing in $k$ and attains its maximum when $k$ goes to zero. At the limit, the Price of Anarchy simplifies to $(e^{\mathcal{E}^{n-1}(1)} - 1)/(\mathcal{E}^{n-1}(1))$.

It remains to show by induction that $\bar{F}(\bar{y}_{i-1}^d) = \bar{F}(\bar{y}_i^d)(1 + g(\bar{y}_{i-1}^d)\mathcal{E}^{i-2}(1))$ for $i = 2, ..., n$ with the proposed piecewise Pareto distribution. To simplify the exposition, we remove the superscript $d$, as all quantities refer to the decentralized case.

Consider supplier 1. Let us fix $\bar{F}(\bar{y}_1) = \varphi$ and $g(\bar{y}_1) = k$. Using Lemma 4, the distribution that minimizes $\int_0^{\bar{y}_1} \bar{F}(x)dx$ is Pareto with a generalized failure rate converging to zero. With that distribution, $\int_0^{\bar{y}_1} \bar{F}(x)dx = \varphi\bar{y}_1$, and the first-order optimality condition for $\bar{y}_1$ simplifies to $\bar{F}(\bar{y}_1) = \bar{F}(\bar{y}_2)(1 + k) = \bar{F}(\bar{y}_2)(1 + k\mathcal{E}^0(1))$.

Let us consider supplier $i \leq n$ and fix $\bar{F}(\bar{y}_i) = \varphi$ and $g(\bar{y}_i) = k$. By induction hypothesis, the distribution between 0 and $\bar{y}_{i-1}$ is piecewise Pareto, with breakpoints

at $\bar{y}_j$, $j \leq i-1$; the generalized failure rates of all pieces, except the last one, tend to zero, and the convergence to zero is faster on the pieces associated with lower values. Accordingly, $\int_0^{\bar{y}_{i-1}} \bar{F}(x)dx = \bar{y}_{i-1}\bar{F}(\bar{y}_{i-1})$, with $\bar{F}(\bar{y}_{i-1}) = \varphi(1 + g(\bar{y}_{i-1})\mathcal{E}^{i-2}(1))$.

Using Lemma 4, the distribution that minimizes $\int_{\bar{y}_{i-1}}^{\bar{y}_i} \bar{F}(x)dx$ when $\bar{F}(\bar{y}_i) = \varphi$, is Pareto between $\bar{y}_{i-1}$ and $\bar{y}_i$ with a generalized failure rate converging to zero, i.e., $\bar{F}_i(x) = (\bar{y}_i/x)^l\varphi$ when $l$ tends to zero. Thus,

$$
\begin{aligned}
\int_{\bar{y}_{i-1}}^{\bar{y}_i} \bar{F}(x)dx &= \frac{\varphi}{1-l}\bar{y}_i(1 - (\frac{\bar{y}_{i-1}}{\bar{y}_i})^{1-l}), \\
&= \frac{\varphi}{1-l}\bar{y}_i(1 - (\frac{\bar{F}(\bar{y}_{i-1})}{\varphi})^{-\frac{1-l}{l}}), \\
&= \frac{\varphi}{1-l}\bar{y}_i(1 - (1 + l\mathcal{E}^{i-2}(1))^{-\frac{1-l}{l}}),
\end{aligned}
$$

where the second equality follows from $\bar{F}(\bar{y}_{i-1}) = (\bar{y}_i/\bar{y}_{i-1})^k\varphi$ by construction, and the next-to-last equality is derived from the induction hypothesis. Therefore

$$
\begin{aligned}
\int_0^{\bar{y}_i} \bar{F}(x)dx &= \int_0^{\bar{y}_{i-1}} \bar{F}(x)dx + \int_{\bar{y}_{i-1}}^{\bar{y}_i} \bar{F}(x)dx, \\
&= \bar{y}_{i-1}\bar{F}(\bar{y}_{i-1}) + \frac{\varphi}{1-l}\bar{y}_i(1 - (1 + l\mathcal{E}^{i-2}(1))^{-\frac{1-l}{l}}), \\
&= \bar{y}_{i-1}\varphi(1 + l\mathcal{E}^{i-2}(1)) + \frac{\varphi}{1-l}\bar{y}_i(1 - (1 + l\mathcal{E}^{i-2}(1))^{-\frac{1-l}{l}}), \\
&= \frac{\varphi}{1-l}\bar{y}_i,
\end{aligned}
$$

by applying twice the induction hypothesis. Therefore, $\int_0^{\bar{y}_i} \bar{F}(x)dx$ tends to $\varphi\bar{y}_i$ as $l \to 0$.

On the other hand, the first-order condition for optimality simplifies to

$$
\begin{aligned}
\bar{F}(\bar{y}_i) &= \bar{F}(\bar{y}_{i+1})\left(1 + \frac{f(\bar{y}_i)}{(\bar{F}(\bar{y}_i))^2}\int_{\bar{y}_{i-1}}^{\bar{y}_i} \bar{F}(x)dx\right), \\
&= \bar{F}(\bar{y}_{i+1})\left(1 + \frac{k}{1-l}(1 - (1 + l\mathcal{E}^{i-2}(1))^{-\frac{1-l}{l}})\right), \\
&= \bar{F}(\bar{y}_{i+1})\left(1 + \frac{k}{1-l}(1 - (1 + l\mathcal{E}^{i-2}(1))^{\frac{-1}{l}}(1 + l\mathcal{E}^{i-2}(1)))\right),
\end{aligned}
$$

Taking the limit when $l$ goes to zero simplifies the first-order condition to

$$\bar{F}(\bar{y}_i) = \bar{F}(\bar{y}_{i+1})\left(1 + k\mathcal{E}^{i-1}(1))\right),$$

since $\lim_{l\to 0}(1 + xl)^{(-1/l)} = e^{-x}$, preserving the induction hypothesis for the next supplier. □

## Proof of Theorem 16

The proof follows the proof of Theorem 10. In Lemma 6, the worst-case distribution is such that $f(x) \to 0$ and $\bar{F}(x) \to \bar{F}(y^d)$, for $x < y^d$. Accordingly, (2.5) reduces to $\bar{F}(y^d)(1 - g(y^d)/n) = r$. Replacing $\varphi$ with $r/(1 - k/n)$ in upper bound (A.2), we obtain an increasing function of $k$. Since $\bar{F}(y^d) \leq 1$, $k \leq (1 - r)n$. Replacing $k$ with its upper bound in the ratio gives rise to the theorem statement. □

## Proof of Theorem 17

The proof is similar to the proof of Theorem 10. Using (2.6), we replace $\varphi$ with $r/(1 - k/n)$ in (A.2) of Lemma 6. The resulting upper bound is increasing with $k$ and maximized when $k = (1 - r)n$ (since $r/(1 - k/n) = \bar{F}(y^d) \leq 1$), leading to the theorem statement. □

## Proof of Theorem 18

The proof is similar to the proof of Theorem 12. In Lemma 6, the worst-case distribution is such that $\int_0^{y^d} \bar{F}(x)dx = \bar{F}(y^d)y^d$. Thus, optimality condition (2.7) simplifies to $\bar{F}(y^d) = r(1 + nk)$, where $k = g(y^d)$. Replacing $\varphi$ with $r(1 + nk)$ into upper bound (A.2) simplifies the Price of Anarchy to $((1 - kn)^{1/k} - (1 - kn))/(n(1 - k))$, which is decreasing with $k$ and is then maximized when $k \to 0$. □

192

# A.3 Chapter 3

## Existence of at most one shock in the case of a quadratic fundamental diagram

Potentially, a shock could result from the focusing of two even forward waves, or two even backward waves. We will impose a condition eliminating these cases. This will imply that there is at most one shock on every road, occurring when forward waves intersect backward waves.

Let us first consider two characteristic lines associated with forward waves, one emanating from the origin 0, and the other emanating from some arbitrary location $y$, with respective positive slopes defined according to Assumption 2. That is,

$$\frac{df(k(0,t))}{dk} = v^{max}(1 - 2k(0,t)/k^{max}), \text{ and}$$
$$\frac{df(k(y,t))}{dk} = v^{max}(1 - 2(k(0,t) + B(0,t)y)/k^{max}).$$

A shock occurs at the intersection of these two characteristic lines, at location $\hat{x} = (k^{max} - 2k(0,t))/(2B(0,t))$. Since the shock location is independent of $y$, we conclude that all forward characteristic lines intersect each other at the same location.

Therefore, no local shock occurs if $\hat{x} > L$, that is, if $B(0,t) \leq (k^{max}/2 - k(0,t))/L$. This is automatically satisfied if $0 \leq k(0,t) + B(0,t)L \leq k^{max}/2$.

Similarly, for backwards moving waves, we consider two characteristic lines, one emanating from the destination $L$, and the other emanating from some location $0 < y < L$. No local shock occurs if the shock location is below 0 (before the road entrance), i.e., $B(L,t) \leq (k(L,t) - k^{max}/2)/L$. This is automatically satisfied if $k^{max}/2 \leq k(L,t) - B(L,t)L \leq k^{max}$. $\square$

## Proof of Theorem 19

The travel time to go from the road entrance to the shock location is $u_0\hat{x}$. On the other hand, the vehicle's trajectory in the heavy traffic region is given by differen-

tial equation (3.11). Solving the differential equation, together with the boundary condition $\tau(\hat{x}, t_0) = u_0 \hat{x}$ and using a power series solution gives rise to (3.12).

This power series is guaranteed to converge as long as Assumption 4 holds. Indeed, the ratio of two successive terms is bounded from above by

$$\left| \frac{2B(L, t_0 + \theta) k^{max}(L - \hat{x})}{(k^{max} - k(L, t_0 + \theta) + B(L, t_0 + \theta)(L - \hat{x} + (\theta - u_0 \hat{x})/w_0))^2} \right|,$$

and the series converges if the ratio is bounded from above by 1.

If $B(L, t_0 + \theta) \geq 0$, an upper bound on the denominator of the ratio is $(k^{max} - k(L, t_0 + \theta))^2$ since both $L - \hat{x}$ and $\theta - u_0 \hat{x}$ are positive, by assumption on $\theta$. Hence, the ratio is bounded from above by 1 if $B(L, t_0 + \theta) < (k^{max} - k(L, t_0 + \theta))^2 / (2Lk^{max})$, which holds under Assumption 4.

If $B(L, t_0 + \theta) \leq 0$, the term $(L - \hat{x} + (\theta - u_0 \hat{x})/w_0)$ is bounded from above by $(L - \hat{x})(k^{max}/(k^{max} - k(L, t_0 + \theta)))$, by assumption on $\theta$. Replacing $B(L, t_0 + \theta)$ by its lower bound (3.6) in the denominator, we obtain that the denominator is bounded from below by $(3/4(k^{max} - k(L, t_0 + \theta))^2$. Hence, the above ratio is less than 1 if $-B(L, t_0 + \theta) < 9/16(k^{max} - k(L, t_0 + \theta))^2 / (2Lk^{max})$, which holds under Assumption 4. $\qquad \square$


## Proof of Theorem 20

We obtain an ODE representing the instantaneous velocity of a vehicle in a similar fashion as in the case of light traffic. However, the boundary condition is defined at the shock location, $\tau(\hat{x}, t_0) = \hat{\tau}$, where $\tau(\hat{x}, t_0)$ is the time to reach the shock, using (3.13). Using a power series solution to solve the differential equation, together with the boundary condition, gives rise to (3.15). The ratio between two successive terms in this power series is bounded from above by

$$\left| \frac{2B(L, t_0 + \theta)(k^{max} + 2v^{max}B(L, t_0 + \theta)(\theta - \hat{\tau}))(L - \hat{x})}{(k^{max} - k(L, t_0 + \theta) + B(L, t_0 + \theta)(L - \hat{x} + v^{max}(\theta - \hat{\tau})))^2} \right|,$$

and the series converges if this ratio is less than 1.

194

If $B(L, t_0 + \theta) \geq 0$, the denominator of the ratio is bounded from above by $(k^{max} - k(L, t_0 + \theta))^2$ since both $L - \hat{x}$ and $\theta - \hat{\tau}$ are positive, by assumption on $\theta$. Replacing $B(L, t_0 + \theta)$ by its upper bound (3.6) in the numerator, we obtain that the numerator is bounded from above by $2B(L, t_0 + \theta)(k^{max} + 1/2(k^{max} - k(L, t_0 + \theta))) < 5/2B(L, t_0 + \theta)k^{max}$, since $k(L, t_0 + \theta) \geq k^{max}/2$. Hence, the ratio is bounded by 1 if $B(L, t_0 + \theta) < 2/5(k^{max} - k(L, t_0 + \theta))^2/(Lk^{max})$, which holds under Assumption 4.

If $B(L, t_0 + \theta) \leq 0$, the term $(L - \hat{x} + v^{max}(\theta - \hat{\tau}))$ is bounded from above by $(L - \hat{x})(2k^{max} - k(L, t_0 + \theta))/(k^{max} - k(L, t_0 + \theta))$, by assumption on $\theta$. Replacing $B(L, t_0 + \theta)$ by its lower bound (3.6) in the denominator, we obtain that the denominator is bounded from below by $((k^{max} - k(L, t_0 + \theta))(1 - 1/4(2 - k(L, t_0 + \theta)/k^{max})))^2$. Since $k(L, t_0 + \theta) \leq k^{max}$, the denominator is bounded from below by $(3/4(k^{max} - k(L, t_0 + \theta)))^2$. Hence, the ratio is less than 1 if $-B(L, t_0 + \theta) < 9/16(k^{max} - k(L, t_0 + \theta))^2/(2Lk^{max})$, which holds under Assumption 4. $\qquad\square$

## Proof of Theorem 21

The shock locations (3.10) and (3.14) are continuous functions of $A(0, t_0)$. Since both travel-time functions in light traffic ($xu_0$ and (3.13)) are continuous functions of $x$ and $k(0, t)$, and the travel-time functions (3.12) and (3.15) are continuous functions of $\hat{\tau}$, it follows that the travel-time functions are continuous. $\qquad\square$

## Proof of Theorem 22

Proving the (strict) monotonicity of a function is equivalent to showing that its Jacobian matrix is positive semi-definite (positive definite) (see Nagurney 1993). The $(i, j)$-th entry of the Jacobian is the derivative of the travel time of the $i$th period with respect to the flow incoming in the $j$th period. In what follows, we consider the derivatives of the travel times with the incoming densities, as there is a one-to-one relation between the incoming flow and the incoming density (recall that, from Assumption 3, we always have $k(0, t) \leq k(f^{max})$).

In the remainder of the proof, we show that, if $B(0, t)$ and $B(L, t + \theta)$ are in the

some intervals around zero, the Jacobian matrix is lower triangular, with nonnegative diagonal entries. In particular, if the fundamental diagram is quadratic, all diagonal entries are positive, and the Jacobian is positive definite. If the fundamental diagram is triangular, some diagonal entries might be equal to zero, and the Jacobian matrix is positive semidefinite.

Let us first consider the case without shocks. The free flow travel time consistent with a triangular fundamental diagram is independent of flow. As a result, if there is no shock in period $i$, the entire $i$-th row in the matrix is zero. On the other hand, with a quadratic fundamental diagram, the travel time under light traffic depends on the entering flow, both through $k(0,t)$ and $B(0,t)$. In the following, we assume that $B(0, t+\gamma(t)) = (k(0,t) - k(0, t+\gamma(t)))/L$, with $\gamma(t) = L/(df(0,t)/dk) \geq 0$. (Simpler results can be derived if $B(0,t)$ is independent of the incoming flow.) The travel time in period $t + \gamma(t)$ is a function of $k(0,t)$ and $k(0, t+\gamma(t))$; hence the travel time depends only on past flows. Considering the first-order travel-time function (3.13), the diagonal elements of the Jacobian matrix are of the form

$$
\frac{d\tau(L,t)}{dk(0,t)} = \frac{L(k^{max} - (3/2)k(0,t))k^{max}}{v^{max}(k^{max} - k(0,t))^3} + \frac{B(0,t)L^2(k^{max} + 2k(0,t))k^{max}}{2v^{max}(k^{max} - k(0,t))^4}.
$$

Therefore, if $B(0,t)$ is in some interval $[-\underline{B}_{low}, \infty)$ around zero, $\frac{d\tau(L,t)}{dk(0,t)} > 0$. The travel times in periods other than $t$ and $t + L/(df(0,t)/dk)$ are independent of $k(0,t)$; hence, their derivative with respect to $k(0,t)$ is zero.

Let us now introduce shocks into the travel-time functions. Independently of the shape of the fundamental diagram, the shock location for a vehicle departing at time $t$, (3.10) and (3.14), decreases linearly with the flows that have started in periods $s \leq t$ and that have not left the road at $t + \theta$, i.e., $s + \tau(L,s) < t + \theta$, through $A(0,t) - D(L, t + \theta)$.

In addition, one can show that the travel time with shocks is a decreasing function of the shock location, if $B(L, t + \theta)$ is in some interval $[-\underline{B}_{high}, \overline{B}_{high}]$ around zero. In the case of a triangular fundamental diagram, the derivative of the first order

travel-time function (3.12) with respect to the shock location,

$$\frac{d\tau(L,t)}{d\hat{x}} = u_0 + w_0 - k^{max}\frac{w_0(k^{max} - k(L,t+\theta)) + B(L,t+\theta)(\theta - u_0 L)}{(k^{max} - k(L,t+\theta) + B(L,t+\theta)(L - \hat{x} + \frac{\theta - u_0\hat{x}}{w_0}))^2},$$

is negative if $B(L,t+\theta)$ is in some interval $[-\underline{B}_{high}, \overline{B}_{high}]$ around zero. If we multiply the first term $u_0 + w_0$ by the denominator, we obtain a quadratic equation in terms of $B(L,t+\theta)$, with positive and negative roots. Similarly, in the case of a quadratic fundamental diagram, the travel time (3.15) is decreasing with the shock location $\hat{x}$, for some range of values of $B(L,t+\theta)$ around zero.

As a result, a flow increase in period $s \le t$, with $s + \tau(L,s) \le t + \theta$, will decrease the shock location, and hence increase the travel time of a vehicle departing at time $t$.

Thus, the proposed travel-time functions only depend on past flows. Moreover, their derivative with respect to the current flow is positive in the case of a quadratic fundamental diagram, and nonnegative in the case of a triangular fundamental diagram. As a result, the Jacobian matrix is positive (semi-)definite. □

## Proof of Theorem 23

Satisfying FIFO means that the derivative of the travel time with respect to time, i.e., $d\tau/dt$ is greater than or equal to -1. For the triangular fundamental diagram, the free flow travel time does not depend on previous flows, and the travel-time function is defined by (3.17). Because $f(0,t)/f(L,t+\theta) > 0$, the difference between $\tau(L,t)$ and $\tau(L,t-1)$ is greater than $-1$. The travel-time function in the case of a quadratic fundamental diagram (3.15) is a polynomial of $B(0,t)$, and so is $d\tau/dk$. Since FIFO is satisfied when $B(0,t) = 0$, as (3.15) simplifies to (3.17), there is a range of values around zero for $B(0,t)$ for which the property is also satisfied. □

# Bibliography

Acemoglu, D. and Ozdaglar, A. (2004). Flow control, routing and performance with a for-profit service provider. Working Paper, MIT.

Anupindi, R. and Bassok, Y. (1999). Centralization of stocks: Retailers vs. manufacturer. *Management Sci.*, 45(2):178–191.

Arrow, K., Harris, T., and Marschak, J. (1951). Optimal inventory policy. *Econometrica*, 19(3):250–272.

Arrow, K. J. (1958). Historical background. In Arrow, K. J., Karlin, S., and Scarf, H. E., editors, *Studies in Mathematical Theory of Inventory and Production*, pages 3–15. Stanford University Press, Stanford, CA.

Averbakh, I. (2001). On the complexity of a class of combinatorial optimization problems with uncertainty. *Mathematical Programming*, 90(2):263–272.

Barlow, R. E. and Proschan, F. (1965). *Mathematical Theory of Reliability.* J. Wiley & Sons, New York.

Bell, D. E. (1982). Regret in decision making under uncertainty. *Oper. Res.*, 30(5):961–981.

Belobaba, P. P. (1987). *Air Travel Demand and Airline Seat Inventory Management.* PhD thesis, Flight Transportation Laboratory, MIT, Cambridge, MA.

Belobaba, P. P. (1992). Optimal vs. heuristic methods for nested seat allocation. In *Proc. AGIFORS Reservations and Yield Management Study Group*, Brussels, Belgium.

Ben-Tal, A. and Nemirovski, A. (1999). Robust solutions of uncertain linear programs. *Oper. Res. Lett.*, 25(1):1–13.

Bergemann, D. and Schlag, K. (2005). Robust monopoly pricing: The case of regret. European University Institute Working Paper.

Bernstein, F. and DeCroix, G. A. (2004). Decentralized pricing and capacity decisions in a multitier system with modular assembly. *Management Sci.*, 50(9):1293–1308.

Bernstein, F. and Federgruen, A. (2005). Decentralized supply chains with competing retailers under demand uncertainty. *Management Sci.*, 51(1):18–29.

Bernstein, F. and Marx, L. M. (2005). Reservation profit levels and the division of supply chain profit. Working Paper, Duke.

Bertrand, J. (1883). Théorie mathématique de la richesse sociale. *Journal des Savants*, 67:499–508.

Bertsekas, D. P. (1999). *Nonlinear Programming*. Athena Scientific, Belmont, MA. 2nd edition.

Bertsimas, D. and de Boer, S. (2005). Simulation-based booking limits for airline revenue management. *Oper. Res.*, 53(1):90–106.

Bertsimas, D. and Popescu, I. (2002). On the relation between option and stock prices: A convex optimization approach. *Oper. Res.*, 50(2):358–374.

Bertsimas, D. and Popescu, I. (2003). Revenue management in a dynamic network environment. *Transp. Sci.*, 37(3):257–277.

Bertsimas, D. and Popescu, I. (2005). Optimal inequalities in probability theory: A convex optimization approach. *SIAM J. of Optimization*, 15(3):780–804.

Bertsimas, D. and Sim, M. (2004). The price of robustness. *Oper. Res.*, 52(1):35–53.

Bertsimas, D. and Thiele, A. (2006). A robust optimization approach to inventory theory. *Oper. Res.*, 54(1):150–168.

Boyd, E. A. and Bilegan, I. C. (2003). Revenue management and e-commerce. *Management Sci.*, 49(10):1363–1386.

Boyd, S. and Vandenberghe, L. (2004). *Convex Optimization*. Cambridge University Press, United Kingdom.

Brown, A. O. and Tang, C. S. (2006). The impact of alternative performance measures on single-period inventory policy. *J. of Industrial and Management Optimization*. To Appear.

Brumelle, S. L. and McGill, S. I. (1993). Airline seat allocation with multiple nested fare classes. *Oper. Res.*, 4(1):127–137.

Bryant, J. (1980). Competitive equilibrium with price setting firms and stochastic demand. *Internat. Econom. Rev.*, 21(3):619–26.

Cachon, G. P. (2003). Supply chain coordination with contracts. In Graves, S. and de Kok, T., editors, *Handbook in Oper. Res. and Management Sci.: Supply Chain Management*, chapter 6. Elsevier, North-Holland.

Cachon, G. P. (2004). Push, pull and advance-purchase discount contracts. *Management Sci.*, 50(2):222–238.

Cachon, G. P. and Lariviere, M. A. (2001). Contracting to assure supply: How to share demand forecasts in a supply chain. *Management Sci.*, 47(5):629–646.

Carey, M., Ge, Y. E., and McCartney, M. (2003). A whole-link travel-time model with desirable properties. *Transp. Sci.*, 37(1):83–96.

Cayford, R., Lin, W. H., and Daganzo, C. F. (1997). The NETCELL simulation package: technical description. Research report UCB-ITS-PRR-97-23, University of California, Berkeley. URL=http://www.ce.berkeley.edu/~daganzo (last visit: November 2003).

Chamberlain, G. (2000). Econometrics and decision theory. *J. of Econometrics*, 95(2):255–283.

Chan, L. M. A. and Simchi-Levi, D. (2005). Decentralized vs. centralized supply chains. Presentation at the MSOM conference, Northwestern University.

Chen, F., Drezner, Z., Ryan, J. K., and Simchi-Levi, D. (2000). Quantifying the bullwhip effect in a simple supply chain: The impact of forecasting, lead times and information. *Management Sci.*, 46(3):436–443.

Chen, F., Federgruen, A., and Zheng, Y. S. (2001). Coordination mechanisms for a distribution system with one supplier and multiple retailers. *Management Sci.*, 47(5):693–708.

Chi, Z. (1995). *Airline Yield Management in a Dynamic Network Environment.* PhD thesis, MIT, Cambridge, MA.

Conde, E. (2004). An improved algorithm for selecting p items with uncertain returns according to the minmax-regret criterion. *Mathematical Programming*, 100(2):345–353.

Cooper, W. L. (2002). Asymptotic behavior of an allocation policy for revenue management. *Oper. Res.*, 50(4):720–727.

Correa, J. R., Schulz, A. S., and Stier-Moses, N. E. (2004). Selfish routing in capacitated networks. *Math. Oper. Res.*, 29(4):961–976.

Correa, J. R., Schulz, A. S., and Stier-Moses, N. E. (2005). On the efficiency of equilibria in congestion games. In *Proc. IPCO XI*, Lecture Notes in Computer Science, pages 167–181.

Curry, R. E. (1990). Optimal airline seat allcoation with fare classes nested by origins and destinations. *Transp. Sci.*, 24(3):193–204.

Daganzo, C. F. (1994). The cell transmission model: A dynamic representation of highway traffic consistent with the hydrodynamic theory. *Transp. Res. B*, 28(4):269–287.

Daganzo, C. F. (1995a). The cell transmission model. Part II: Network traffic. *Transp. Res. B*, 29(2):79–93.

Daganzo, C. F. (1995b). A finite difference approximation of the kinematic wave model of traffic flow. *Transp. Res. B*, 29(4):261–276.

Daganzo, C. F. (1995c). Properties of link travel time function under dynamic loads. *Transp. Res. B*, 29(2):95–98.

Darlymple, D. J. (1988). Sales forecasting practices. *Internat. J. of Forecasting*, 3(3):379–391.

Dasgupta, P. and Maskin, E. (1986). The existence of equilibrium in discontinuous economic games, I: Theory; II: Applications. *Rev. of Econom. Studies*, 53(1):1–26, 27–41.

de Boer, S. V., Freling, R., and Piersma, N. (2002). Mathematical programming for network revenue management revisited. *Eur. J. Oper. Res.*, 137(1):72–92.

Dell, M. and Magretta, J. (1998). The power of virtual integration: An interview with Dell Computer's Michael Dell. *Harvard Bus. Rev.*, 76(2):72–84.

Ding, X., Puterman, M. L., and Bisi, A. (2002). The censored newsvendor and the optimal acquisition of information. *Oper. Res.*, 50(3):517–527.

Dror, M., Trudeau, P., and Ladany, S. P. (1988). Network models for seat allocation on flights. *Transp. Res. B*, 22(4):293–250.

Edgeworth, F. (1888). The mathematical theory of banking. *J. Royal Statist. Soc.*, 51(1):113–127.

Edwards, Jr., C. H. and Penney, D. E. (1985). *Elementary Differential Equations with Boundary Value Problems*. Prentice Hall, 3rd edition.

Eppen, G. D. and Martin, R. K. (1988). Determining safety stock in the presence of stochastic lead time and demand. *Management Sci.*, 34(11):1380–1390.

Ertogral, K. and Wu, S. D. (2001). A bargaining game for supply chain contracting. Lehigh Working Paper.

Fisher, M. and Raman, A. (1996). Reducing the cost of uncertainty through accurate response to early sales. *Oper. Res.*, 44(1):87–99.

Friesz, T. L., Bernstein, D., Smith, T. E., Tobin, R. L., and Wie, B.-W. (1993). A variational inequality formulation for the dynamic user equilibrium problem. *Oper. Res.*, 41(1):179–191.

Gallego, G. (1992). A minimax distribution-free for the (Q,R) inventory model. *Oper. Res. Lett.*, 11:55–60.

Gallego, G. (1998). New bounds and heuristics for (Q,R) policies. *Management Sci.*, 44(2):219–233.

Gallego, G., Katircioglu, K., and Ramachandran, B. (2006). Inventory management under highly uncertain demand. *Oper. Res. Lett.* To Appear.

Gallego, G. and Moon, I. (1993). The distribution free newsboy problem: Review and extensions. *The J. of the Operational Res. Soc.*, 44(8):825–834.

Gallego, G., Ryan, J. K., and Simchi-Levi, D. (2001). Minimax analysis for finite-horizon inventory models. *IIE Transactions*, 33:861–874.

Gazis, D. and Herman, R. Rothery, R. W. (1961). Nonlinear follow-the-leader models of traffic flow. *Oper. Res.*, 9(4):545–567.

Gerchak, Y. and Wang, Y. (2004). Revenue-sharing vs. wholesale-price contracts in assembly systems with random demand. *Production and Oper. Management*, 13(1):23–33.

Godfrey, G. and Powell, W. B. (2001). An adaptive, distribution-free algorithm for the newsvendor problem with censored demands, with application to inventory and distribution problems. *Management Sci.*, 47:1101–1112.

Graves, S. C. and Parsons, J. C. W. (2005). Using a newsvendor model for inventory planning of NFL replica jerseys. In *Proc. MSOM Conference*, Northwestern University, Evanston, IL.

Greenshields, B. (1935). A study of traffic capacity. In *Highway Research Board Proc.*, volume 14, pages 468–477.

Guo, P. (2003). Newsvendor problems based on possibility theory. In *Knowledge-Based Intelligent Information and Engineering Systems, 7th International Conference*, pages 213–219.

Haberman, R. (1977). *Mathematical Models; Mechanical Vibrations, Population Dynamics and Traffic Flow*. Prentice-Hall.

Hadley, G. and Whitin, T. M. (1963). *Analysis of Inventory Models*. Prentice-Hall, Englewood Cliffs, N.J.

Hayes, R. H. (1969). Statistical estimation problems in inventory control. *Management Sci.*, 15(11):686–701.

Herman, R., Montroll, E. W., Potts, R. B., and Rothery, R. W. (1959). Traffic dynamics: Analysis of stability in car following. *Oper. Res.*, 7:86–106.

Hodges, S. D. and Moore, P. G. (1970). The product-mix problem under stochastic seasonal demand. *Management Sci.*, 17(2):B107–B114.

Hurdle, V. F. and Son, B. (2000). Road test of a free way model. *Transp. Res. A*, 34(7):537–564.

Huyett, W. I. (2005). Extreme competition. *McKinsey Quarterly*, (1):46–57.

Inuiguchi, M. and Sakawa, M. (1995). Minimax regret solutions to linear programming problems with an interval objective function. *Eur. J. Oper. Res.*, 86(3):526–536.

Jaynes, E. T. (1957). Information theory and statistical mechanics. *Physical Rev.*, 106:620–630.

Jaynes, E. T. (2003). *Probability Theory: The Logic of Science*. Cambridge University Press.

Johari, R. and Tsitsiklis, J. N. (2004). Network resource allocation and a congestion game. *Math. Oper. Res.*, 29(3):407–435.

Kachani, S. and Perakis, G. (2001). Modeling travel times in dynamic transportation networks; a fluid dynamics approach. Working Paper, MIT.

Karp, R. M. (1992). On-line algorithms vs. off-line algorithms: How much is it worth knowing the future? In van Leeuwen, J., editor, *Proc. World Computer Congress*, volume 1, pages 416–429. Elsevier Science Publishers.

Kasugai, H. and Kasegai, T. (1961). Note on minimax regret ordering policy - static and dynamic solutions and a comparison to maximin policy. *J. of the Oper. Res. Soc. of Japan*, 3:155–169.

Keser, C. and Paleologlo, G. A. (2004). Experimental investigation of supplier-retailer contracts: The wholesale price contract. Technical Report 2004s-57, CIRANO, Montreal, QC, Canada.

Keynes, J. M. (1921). *A Treatise on Probability*. MacMillan, New York.

Khoo, B. C., Lin, G. C., Peraire, J., and Perakis, G. (2002). A dynamic user-equilibrium model with travel times computed from simulation. Working Paper, MIT.

Koutsoupias, E. and Papadimitriou, C. H. (1999). Worst-case equilibria. In *16th Symp. on Theoretical Aspects of Computer Science*, pages 404–413.

Kouvelis, P. and Yu, G. (1997). *Robust Discrete Optimization and Applications*. Kluwer Academic Publishers, Boston.

Krishnan, H., Kapuscinski, R., and Butz, D. A. (2004). Coordinating contracts for decentralized supply chains with retailer promotional effort. *Management Sci.*, 50(1):48–63.

Kuwahara, M. and Akamatsu, T. (2001). Dynamic user optimal assignment with physical queues for a many-to-many OD pattern. *Transp. Res. B*, 35:461–479.

Lancaster, J. (2003). The financial risk of airline revenue management. *J. of Revenue and Pricing Management*, 2(2):158–165.

Lariviere, M. and Porteus, E. (2001). Selling to the newsvendor: An analysis of price-only contracts. *Manufacturing and Service Oper. Management*, 3(4):293–305.

Lariviere, M. A. (1999). Supply chain contracting and coordination with stochastic demand. In Tayur, S., Magazine, M., and Ganeshan, R., editors, *Quantitative Models for Supply Chain Management*, chapter 8, pages 233–268. Kluwer, Doordrecht, North-Holland.

Lariviere, M. A. (2005). A note on probability distributions with increasing generalized failure rates. *Oper. Res.* To Appear.

Lariviere, M. A. and Porteus, E. L. (1999). Stalking information: Bayesian inventory management with unobserved lost sales. *Management Sci.*, 45(3):346–363.

Lau, H.-S. and Lau, A. H.-L. (1995). The multi-product multi-constraint newsboy problem: Application, formulation and solution. *J. of Oper. Management*, 13(3):153–162.

Lau, H.-S. and Lau, A. H.-L. (1996). The newsstand problem: A capacitated multiple-product single-period inventory problem. *Eur. J. of Oper. Res.*, 94(1):29–42.

Levi, R., Roundy, R., and Schmoys, D. B. (2005). Provably near-optimal sample-based policies for stochastic inventory control models. Cornell Working Paper.

Li, J., Fujiwara, O., and Kawakami, S. (2000). A reactive dynamic user equilibrium model in network with queues. *Transp. Res. B*, 34:605–624.

Lighthill, M. J. and Whitham, G. B. (1955). On kinematic waves: II. A theory of traffic flow on long crowded roads. In *Proc. Royal Soc. A*, volume 229, pages 281–345, London.

Lim, A. E. B. and Shanthikumar, J. G. (2006). Relative entropy, exponential utility, and robust dynamic pricing. UC Berkeley Working Paper.

Lin, W. H. and Lo, H. K. (2000). Are the objectives and solutions of dynamic user-equilibrium models always consistent? *Transp. Res. A*, 34:137–144.

Lippman, S. A. and McCardle, K. F. (1997). The competitive newsboy. *Oper. Res.*, 45(1):54–65.

Littlewood, K. (1972). Forecasting and control of passenger bookings. In *Proc. of the Twelfth Annual AGIFORS Symposium*, Nathanya, Israel.

Loomes, G. and Sugden, R. (1982). Regret theory: An alternative theory of rational choice under uncertainty. *The Econom. J.*, 92:805–824.

Luce, R. F. and Raiffa, H. (1957). *Games and Decisions*. John Wiley and Sons, New York.

Mahut, M. (2000). *Discrete Flow Model for Dynamic Network Loading*. PhD thesis, Université de Montréal - Département d'informatique et de recherche opérationnelle.

Martínez-de Albéniz, V. and Simchi-Levi, D. (2003). Competition in the supply option market. Working Paper, MIT.

Martos, B. (1975). *Nonlinear Programming. Theory and Methods*. North Holland Publishing Company, Amsterdam.

Moon, I. and Choi, S. (1995). The distribution free newsboy problem with balking. *The J. of the Operational Res. Soc.*, 46(4):537–542.

Moon, I. and Gallego, G. (1994). Distribution free procedures for some inventory models. *The J. of the Operational Res. Soc.*, 45(6):651–658.

Morris, W. T. (1959). Inventorying for unknown demand. *The J. of Industrial Engineering*, X(4):299–302.

Naddor, E. (1978). Sensitivity to distributions in inventory systems. *Management Sci.*, 24(16):1769–1772.

Nagurney, A. (1993). Network economics: A variational inequality approach. In *Advances in Computational Economics*, volume 10. Kluwer Economic Publishers.

Nahmias, S. and Schmidt, C. P. (1984). An efficient heuristic for the multi-item newsboy problem with a single constraint. *Naval Res. Logist. Quart.*, 31(3):463–474.

Newell, G. F. (1993). A simplified theory of kinematic waves in highway traffic, I. General theory; II Q.ueuing at freeway bottlenecks; III. Multidestination flows. *Transp. Res. B*, 27(4):281–314.

Ouyang, L.-Y. and Wu, K.-S. (1998). A minimax distribution free procedure for mixed inventory model with variable lead-time. *Internat. J. of Production Econom.*, 56-57:511–516.

Papadimitriou, C. H. (2001). Algorithms, games, and the Internet. In *Proc. of the 33rd Annual ACM Symposium on Theory of Cumputing (STOC)*, pages 749–753, Hersonissos, Greece.

Patriksson, M. (1994). The traffic assignment problem: Models and methods. In *Topics in Transportation*. VSP BV, Zeist, The Netherlands.

Paul, A. (2005). A note on closure properties of failure rate distributions. *Oper. Res.*, 53(4):733–734.

Perakis, G. (2000). Dynamic traffic flow problems; a hydrodynamic theory approach. Working Paper, MIT.

Perakis, G. (2005). The "Price of Anarchy" under nonlinear and asymmetric costs. *Math. Oper. Res.* To Appear.

Petrović, D., Petrović, R., and Vujošević, M. (1996). Fuzzy models for the newsboy problem. *Internat. J. of Production Econom.*, 45:435–441.

Popescu, I. (2005). A semidefinite approach to optimal moment bounds for convex class of distributions. *Mathematics of Oper. Res.*, 50. To appear.

Porteus, E. L. (1990). Stochastic inventory theory. In Heyman, D. and Sobel, M., editors, *Handbook in Oper. Res. and Management Sci.*, volume 2, pages 605–653. Elsevier North-Holland, Amsterdam.

Porteus, E. L. (2002). *Foundations of Stochastic Inventory Management*. Stanford University Press, Stanford, CA.

Ran, B. and Boyce, D. E. (1994). Dynamic urban transportation network models. In *Lecture Notes in Economics and Mathematical Systems*, volume 417. Springer-Verlag.

Ran, B., Rouphail, N. M., Tarko, A., and Boyce, D. E. (1997). Toward a class of link travel time functions for dynamic assignment models on signalized networks. *Transp. Res. B*, 31(4):277–290.

Richards, P. I. (1956). Shock waves on the highway. *Oper. Res.*, 4:42–51.

Robinson, L. W. (1995). Optimal and approximate control policies for airline booking with sequential nonmonotoic fare classes. *Oper. Res.*, 43(2):252–263.

Roughgarden, T. and Tardos, E. (2000). How bad is selfish routing? In *41st IEEE Symp. on Foundations of Computer Science*, pages 93–102.

Roughgarden, T. and Tardos, E. (2002). Bounding the inefficiency of equilibria in nonatormic congestion games. *J. of the ACM*, 49(2):236–259.

Savage, L. J. (1951). The theory of statistical decisions. *J. of the Amer. Stat. Assoc.*, 46(253):55–67.

Savage, L. J. (1954). *The Foundations of Statistics*. J. Wiley and Sons, New York.

Scarf, H. (1960). Some remarks on Bayes solutions to the inventory problem. *Naval Res. Logist. Quart.*, 7:591–596.

Scarf, H. E. (1958). A min-max solution to an inventory problem. In Arrow, K. J., Karlin, S., and Scarf, H. E., editors, *Studies in Mathematical Theory of Inventory and Production*, pages 201–209. Stanford University Press, Stanford, CA.

Schrank, D. and Lomas, T. (2003). The 2003 annual urban mobility report. Technical report, Texas Transportation Institute.

Schweitzer, M. A. and Cachon, G. P. (2000). Decision bias in the newsvendor problem with a known demand distribution: Experimental evidence. *Management Sci.*, 46(3):404–420.

Simpson, R. W. (1989). Using network flow techniques to find shadow prices for market and seat inventory control. Technical Report M89-1, MIT Flight Transportation Laboratory, Cambridge, MA.

Smith, J. E. (1995). Generalized Chebyshev inequalities: Theory and applications in decision analysis. *Oper. Res.*, 43(5):807–825.

Spengler, J. (1950). Vertical integration and antitrust policy. *J. of Political Economy*, 4(58):347–352.

Talluri, K. I. and van Ryzin, G. J. (1998). An analysis of bid-price controls for network revenue management. *Management Sci.*, 44(11):1577–1593.

Talluri, K. I. and van Ryzin, G. J. (1999). A randomized linear programming formulation method for computing network bid prices. *Transp. Sci.*, 33(2):207–216.

Talluri, K. T. and van Ryzin, G. J. (2004). *The Theory and Practice of Revenue Management*. Kluwer Academic Publishers, Boston, MA.

Tomlin, B. (2003). Capacity investment in supply chains: Sharing the gain rather than sharing the pain. *Manufacturing Service Oper. Management*, 5(4):317–333.

Vairaktarakis, G. L. (2000). Robust multi-item newsboy models with a budget constraint. *Internat. J. of Production Econom.*, 66:213–226.

Van Mieghem, J. A. and Rudi, N. (2002). Newsvendor networks: Inventory manage-mement and capacity investement with discretionary activities. *Manufacturing and Service Oper. Management*, 4(4):313–335.

van Ryzin, G. J. and Vulcano, G. (2005). Simulation-based optimization of virtual nesting controls for network revenue management. NYU Working Paper.

Velan, S. (2000). *The Cell Transmission Model: A New Look at a Dynamic Network Loading Model*. PhD thesis, CRT, Université de Montréal.

Wang, Q. (2001). Coordinating independent buyers in a distribution system to in-crease a vendor's profits. *Manufacturing Service Oper. Management*, 3(4):337–348.

Wang, Y. and Gerchak, Y. (2003). Capacity games in assembly systems with uncertain demand. *Manufacturing Service Oper. Management*, 5(3):252–267.

Wardrop, J. G. (1952). Some theoretical aspects of road traffic research. In *Proc. Inst. Civ. Eng., Part II*, volume 1, pages 325–378.

Williamson, E. L. (1988). Comparison of optimization techniques for origin-destination seat inventory comtrol. Master's thesis, MIT, Cambridge, MA.

Williamson, E. L. (1992). *Airline Network Seat Control: Methodologies and Revenue Impacts*. PhD thesis, MIT, Cambridge, MA.

Wollmer, R. D. (1986). A hub-and-spoke seat management model. Technical report, Douglas Aircraft Company, McDonnell Douglas Corporation.

Wollmer, R. D. (1992). An airline seat management model for a single leg route when lower fare classes book first. *Oper. Res.*, 40(1):26–37.

Yu, G. (1997). Robust economic order quantity models. *Eur. J. Oper. Res.*, 100:482–493.

Yue, J., Chen, B., and Wang, M.-C. (2006). Expected value of distribution informa-tion for the newsvendor problem. *Oper. Res.* To Appear.

Zhou, L. and Natarajan, K. (2005). Robust single period newsvendor model. Technical report, National University of Singapore, Singapore. Undergraduate Research Opportunity Program in Science.

Ziliaskopoulos, A. K. (2000). A linear programming model for the single destination system optimum dynamic traffic assignment problem. *Transp. Sci.*, 34(1):37–49.