# UCLA
## UCLA Electronic Theses and Dissertations

**Title**
Strategic Pricing and Resource Allocation: Framework and Applications

**Permalink**
https://escholarship.org/uc/item/5rw625xc

**Author**
Ren, Shaolei

**Publication Date**
2012

Peer reviewed|Thesis/dissertation

UNIVERSITY OF CALIFORNIA

Los Angeles

# Strategic Pricing and Resource Allocation: Framework and Applications

A dissertation submitted in partial satisfaction

of the requirements for the degree

Doctor of Philosophy in Electrical Engineering

by

**Shaolei Ren**

2012

# Strategic Pricing and Resource Allocation: Framework and Applications

by

## Shaolei Ren

Doctor of Philosophy in Electrical Engineering

University of California, Los Angeles, 2012

Professor Mihaela van der Schaar, Chair

Enabled by ubiquitous broadband connectivity and seamless wireless connections, we have witnessed in the past few years the emergence of a plethora of wireless applications, ranging from data communications and social networking to the more recently wireless cloud computing. The growing tension between the exploding demand for such wireless applications and the increasingly scarce network resources (e.g., spectrum, power) has urged a rethinking of the service providers' pricing strategies and network resource management techniques to cope with potential threats of quality-of-service degradation and revenue decreases. Specifically, it has become of paramount importance for service providers to strategically redesign their pricing policies and to understand how various pricing policies will affect the service demand, competition in the market, as well as the network resource management.

In this dissertation, I propose a novel framework to optimize a service provider's pricing policy as well as its network resource allocation decision for profit maximization, in the presence of self-interested participating users that strategically respond to the charged price to maximize their own benefits. Applicable to both static and stochastic environments, the proposed framework explicitly takes into account user

heterogeneity, which is observed in a wide range of applications. Based on the framework, I investigate the problem of optimizing pricing and resource allocation for the service provider's profit maximization in various contexts, including cooperative relay networks, communications markets, online user-generated content platforms, and mobile cloud computing systems.

The dissertation of Shaolei Ren is approved.

Ali H. Sayed

Jason L. Speyer

William Zame

Philip A. Chou

Mihaela van der Schaar, Committee Chair

University of California, Los Angeles

2012

*To my wife.*

TABLE OF CONTENTS

## LIST OF FIGURES

# LIST OF TABLES

# ACKNOWLEDGMENTS

Numerous people have supported and helped me during the development of this dissertation. A few words here are far from enough to express all my appreciation.

First of all, I would like to thank my advisor, Professor Mihaela van der Schaar, for her generous support throughout my study at UCLA and guiding me to learn the art of research and the philosophy behind it. I am much indebted for her patience, criticism and encouragement. Her suggestions and advices on my research are very greatly appreciated. I would also like to thank the other members of my dissertation committee, Professor Ali H. Sayed, Professor Jason L. Speyer, Professor William Zame, and Dr. Philip A. Chou for their time and efforts in evaluating my work.

Secondly, I would be remiss not to express my sincere gratitude to all of the amazing colleagues I have met in the UCLA Networks, Economics, Communication Systems, Informatics and Multimedia Research Lab. I would especially like to thank Dr. Hyunggon Park, Dr. Hsien-Po Shiang, Dr. Yi Su, Dr. Fangwen Fu, Dr. Jaeok Park, Dr. Nicholas Mastronarde, Mr. Yu Zhang, Mr. Yuanzhang Xiao, Mr. Jie Xu, Mr. Jianyu Wang, and Mr. Siming Song whose friendship and support have always been invaluable. I have learnt a lot from them, through their valuable discussions and enlightenments, and have spent four exciting and enjoyable years with them. I would also like to thank my mentors Dr. Yuxiong He and Dr. Sameh Elnikety at Microsoft Research, Redmond, for the opportunity to do exciting work during my summer internship.

Finally, I wish to express my deepest gratitude to my family for their love, encouragement and support throughout my graduate studies. I have been truly blessed and am extremely grateful to have the family. Most importantly, I am incredibly grateful for the love and support of my wife, Jingfei Zhang, with whom I shared countless days full of joy and happiness in Los Angeles. She stayed up many late nights with

me while I worked, and was always there when I needed her most. Her undying love, companionship, and belief in me have given me the power to conquer everything. She is truly the heaven's gift in my life.

<div align="center">VITA</div>

| | |
|---|---|
| 09/2002–07/2006 | **Tsinghua University** |
| | B.E. in Electronic Engineering |
| 09/2006–08/2008 | **The Hong Kong University of Science and Technology** |
| | M.Phil. in Electronic and Computer Engineering |
| 06/2011–09/2011 | **Microsoft Research** |
| | Research Intern in eXtreme Computing Group |
| 09/2008–06/2012 | **University of California, Los Angeles** |
| | Research Assistant in Electrical Engineering |

<div align="center">PUBLICATIONS</div>

Shaolei Ren, J. Park, and M. van der Schaar, "Entry and Spectrum Sharing Scheme Selection in Femtocell Communications Markets," *IEEE/ACM Transactions on Networking*, to appear.

Shaolei Ren and M. van der Schaar, "Data Demand Dynamics in Communications Markets," *IEEE Transactions on Signal Processing*, vol. 60, no. 4, pp. 1986-2000, Apr. 2012.

Shaolei Ren and M. van der Schaar, "Pricing and Distributed Power Control in Wireless Relay Networks," *IEEE Transactions on Signal Processing*, vol. 59, no. 6, pp. 2913-2926, June 2011.

Shaolei Ren and M. van der Schaar, "Distributed Power Allocation in Multi-User Multi-Channel Cellular Relay Networks," *IEEE Transactions on Wireless Communications*, vol. 9, no. 6, pp. 1952-1964, Jun. 2010.

Shaolei Ren and M. van der Schaar, "Pricing and Investment for Online TV Content Platforms," *IEEE Transactions on Multimedia*, submitted.

Shaolei Ren, J. Park, and M. van der Schaar, "Maximizing Profit on User-Generated Content Platforms With Participant Heterogeneity," *IEEE Infocom 2012*.

Shaolei Ren and M. van der Schaar, "Revenue Maximization in a Customer-to-Customer Market with Customer Rationality," *GameNets 2012*.

Shaolei Ren, Y. He, and F. Xu, "Provably-Efficient Job Scheduling for Energy and Fairness in Geographically Distributed Data Centers," *ICDCS 2012*.

Shaolei Ren and M. van der Schaar, "Impacts of Congestion on Data Demand Dynamics in Communications Markets," *IEEE CDC 2011*.

Shaolei Ren, J. Park, and M. van der Schaar, "Profit Maximization on Online Content Platforms," *Allerton Conf. 2011*.

Shaolei Ren, J. Park, and M. van der Schaar, "User Subscription Dynamics and Revenue Maximization in Communication Markets," *IEEE Infocom 2011*.

Shaolei Ren, F. Fu, and M. van der Schaar, "Traffic-Dependent Pricing for Delay-Sensitive Multimedia Networks," *IEEE Globecom 2011*.

Shaolei Ren, J. Park, and M. van der Schaar, "Subscription Dynamics and Competition in Communication Markets," *ACM NetEcon 2010*.

Shaolei Ren, J. Park, and M. van der Schaar, "Dynamics of Service Provider Selection in Communication Markets," *IEEE Globecom 2010*.

Shaolei Ren and M. van der Schaar, "Pricing and Distributed Power Control for Relay Networks," *IEEE ICC 2010*.

Shaolei Ren and M. van der Schaar, "Revenue Maximization and Distributed Power Allocation in Cognitive Radio Networks," *ACM CoRoNet co-located with Mobicom 2009 (invited)*.

# CHAPTER 1

# Introduction

## 1.1 Motivation

Due to the availability of ubiquitous broadband connectivity and seamless wireless connections, we have witnessed in the past few years the emergence of a plethora of wireless applications, ranging from data services and social networking to the more recently wireless cloud computing, which leverage the power of the Internet and wireless networks to facilitate social interactions, information and knowledge sharing, as well as business activities. Notable examples include AT&T Cloud, Google Android, Apple iOS, Facebook and many others. In order to fully realize the benefits of the current system infrastructure, tremendous technological efforts have been dedicated in the past decade to increasing the service provider's profit and to enhancing the system performance in terms of throughput, response time, reliability and many other factors. Among the efforts, resource allocation has received a considerable amount of attention due to its effectiveness in significantly boosting the system performance and hence attracting more users as well as increasing the profitability of the system. The traditional resource allocation approach largely relies on system-wide centralized management, which requires all the users to cooperatively follow the prescribed mechanism or protocol. Nevertheless, in the absence of a central controller, rational or selfish users have incentives to optimize their own performances independently, without considering the social welfare and thus, the existing centralized approaches are no longer applicable in such settings. Moreover, treating pricing-based demand-side management and the

system resource management separately, the existing research does not reap the potential benefits of using pricing as a lever to enable more efficient resource management and profit maximization. In this dissertation, we aim to formalize a new design framework that weaves together strategic pricing and resource allocation, which we shall show can significantly improve the system performance (e.g., in terms of profitability) compared to the state-of-the-art approaches.

## 1.2   Key Challenges

In what follows, we state the key design challenges that are common to various applications.

- **User decision dynamics:** A majority of the existing research on pricing assumes that the *equilibrium* demand is directly achieved as a function of the charged price without considering the underlying dynamics that is required prior to reaching the equilibrium demand state. Essentially, without having the complete information regarding each other, self-interested users engage in a non-cooperative game and behave strategically. Thus, the existence of a (Nash) equilibrium and whether the users' non-cooperative behaviors will lead to an equilibrium are particularly important in this regard. Nevertheless, taking into account the users' dynamic behaviors will couple the service provider's pricing decision with the users' self-interested decisions, thereby significantly complicating the formal analysis and design.

- **User heterogeneity:** What adds to the design challenges is user heterogeneity, which is often neglected in the existing research pertaining to pricing-based resource allocation. Heterogeneity is widely observed in various applications and in different senses: wireless transmitters have heterogeneous channel conditions in relay networks; wireless subscribers have heterogeneous data service demand in communications markets; and content producers have heterogeneous content quality on user-generated con-

tent platforms. The prevailing presence of user heterogeneity makes it difficult to consolidate all the users into one "super user" as well as to design an appropriate pricing scheme to achieve the global optimality.

- **Random environment with unknown dynamics:** Achieving the global optimality of pricing and system resource management hinges on the availability of future information as well as the stochastic knowledge of the underlying dynamic environment, which is not readily available in practice. Without such information, it is an intrinsically difficult problem to optimize the long-term system performance. The existing solutions either reply on the prediction of the future information or assume that the dynamic environment in which algorithms will be applied follows a certain stochastic distribution. Nevertheless, predicting the future is typically vulnerable to prediction errors, whereas assuming that the environment follows a certain stochastic distribution cannot adapt to the real environment which may be arbitrarily random.

In summary, a new design framework for optimally setting prices and managing the system resource is necessary, which explicitly takes into account the users' self-interested and competing behaviors, user heterogeneity, and incomplete information about the (possibly random) environment.

## 1.3   Contributions of the Dissertation

In this dissertation, we propose a formal design framework that leverages the power of pricing to enable more efficient resource management. Building upon the Stackelberg leadership model, the design framework explicitly takes into account the strategic behaviors of self-interested users and finds its usage for a wide range of applications, such as cooperative relay networks, wireless communications markets, online user-generated content platforms and mobile cloud computing systems. In particular, instead of directly assuming a demand function of the price, we also explore

and explicitly consider the dynamic process in which self-interested users strategically interact with each other and respond to the charged price. By doing so, the users' self-interested behaviors are aligned with the system designer's goal. Furthermore, by considering user heterogeneity, we use pricing to proactively reshape the users' behaviors/demands and adapt them to the system resource management. Our proposed design framework is applicable to an arbitrarily random environment and, without the necessity of predicting the future information, it is provably efficient in the sense that the resulting performance loss with respect to the optimal offline algorithm with perfect future information is upper bounded.

In the following, we summarize the remainder of this dissertation that instantiates the proposed design framework using various applications.

### 1.3.1 Chapter 2: Pricing and Power Control in Wireless Relay Networks

In Chapter 2 of the dissertation, we consider a wireless amplify-and-forward relay network with one relay node and multiple source-destination pairs/users and propose a pricing framework that enables the relay to set prices to maximize either its revenue or any desirable system utility. Specifically, depending on the quality of the received signals, the relay sets prices and correspondingly charges the users utilizing its resources for their transmissions. The price is determined in such a way that the relay's revenue or system utility is maximized. Given the specified price, the users competitively employ the relay node to forward their signals. We model each user as a rational player, which aims at maximizing its own net utility through power allocation, and analyze the competition among the users within the framework of non-cooperative game theory. It is shown that, in the game played by the users, there always exists a unique pure Nash equilibrium point that can be achieved through distributed iterations. Next, subject to the availability of complete information about the users at the relay, we propose a low-complexity uniform pricing algorithm and an optimal differentiated pricing algorithm,

in which the relay either charges the users at a sub-optimal uniform price or it charges different prices per user. We also show that, by applying the differentiated pricing algorithm that enforces the users to transmit at certain power levels, any system utility can be maximized. Extensive simulations are conducted to quantify the performance of the proposed methods.

### 1.3.2   Chapter 3: Pricing in Wireless Communications Markets

In Chapter 3, we focus on the users' aggregate data demand dynamics in a wireless communications market served by a monopolistic wireless service provider (WSP). Based on the equilibrium data demand, we optimize the WSP's data plans and long-term network capacity decisions to maximize its profit. First, by considering a market where only one data plan is offered, we show that there exists a unique equilibrium in the data demand dynamics regardless of the data plans, and that the convergence of data demand dynamics is subject to the network congestion cost, which is closely related to the WSP's long-term capacity decision. A sufficient condition on the network congestion cost indicates that the WSP needs to provide a sufficiently large network capacity to guarantee the convergence of data demand dynamics. We also propose a heuristic algorithm that progressively optimizes the WSP's data plan to maximize its equilibrium revenue. Next, we turn to a market where two different data plans are offered. It is shown that the existence of a unique equilibrium data demand depends on the data plans, and the convergence of data demand dynamics is still subject to the network congestion cost (and hence, the WSP's network capacity, too). We formalize the problem of optimizing the WSP's data plans and network capacities to maximize its profit. Finally, we discuss the scenario in which the data plans are offered by two competing WSPs and conduct extensive simulations to validate our analysis.

### 1.3.3   Chapter 4: Pricing in Online User-Generated Content Platforms

In Chapter 4, we focus on user-generated content, such as blogs and self-made videos, which has becoming a key component in emerging social media. Specifically, we consider a user-generated content platform monetized through advertising and managed by an intermediary. To maximize the intermediary's profit given rational participants (i.e., content producers and content viewers), we propose a payment scheme in which the intermediary either taxes or subsidizes a content producer an amount of money proportional to the number of views of the producer's content. First, we use a model with a representative content viewer to determine how the content viewers' attention is allocated across available content by solving a utility maximization problem. Then, by modeling the content producers as self-interested agents making independent production decisions, we show that there exists a unique equilibrium in the content production stage, and propose a best-response dynamics to model the decision-making process and to reach the equilibrium. Next, we study the optimal payment scheme (i.e., the payment scheme maximizing the intermediary's profit) that the intermediary chooses taking into account the decisions made by the representative content viewer and the content producers. In particular, by considering the well-known quality-adjusted Dixit-Stiglitz utility function for the representative content viewer, we derive explicitly the optimal payment per content view and characterize analytical conditions under which the intermediary should tax or subsidize the content producers. Finally, we generalize the analysis by considering heterogeneity in terms of production costs among the content producers.

### 1.3.4   Chapter 5: Pricing and Scheduling in Wireless Cloud Computing

In Chapter 5, we consider a wireless cloud computing system in which a profit-maximizing wireless service provider operates a data center and can provide cloud computing

services to its subscribers. In particular, we focus on batch services, which, due to their non-urgent nature, allow more scheduling flexibility than their interactive counterparts. Unlike the existing research that studied *separately* demand-side management and energy cost saving techniques for the wireless cloud (both of which are critical to profit maximization), we propose a provably-efficient Dynamic Scheduling and Pricing (Dyn-SP) algorithm which, using the pricing mechanism as a lever, *proactively* adapts the service demand to workload scheduling in the data center and opportunistically utilizes low electricity prices to process batch jobs for energy cost saving. Without the necessity of predicting the future information (as assumed by some prior works), Dyn-SP can be applied to an arbitrarily random environment in which the electricity price, available renewable energy supply, wireless network capacities provided by base stations may evolve over time as an arbitrary stochastic process. It is proved that, compared to the optimal offline algorithm with future information, Dyn-SP can produce a close-to-optimal long-term profit while bounding the job queue length in the data center. We perform a simulation study based on both traces and randomly generated data to demonstrate the effectiveness of Dyn-SP. In particular, we show both analytically and numerically that a desired tradeoff between the profit and queueing delay can be obtained by appropriately tuning the control parameter. Our results also indicate that, compared to the other algorithms which neglect demand-side management, cooling system energy consumption, or the queue length information, Dyn-SP achieves a higher average profit while incurring (almost) the same average queueing delay.

### 1.3.5    Chapter 6: Conclusion

Chapter 6 concludes the dissertation and includes a discussion about future research directions.

# CHAPTER 2

# Pricing and Distributed Power Control in Wireless Relay Networks

In this chapter, we consider a wireless amplify-and-forward relay network with one relay node and multiple source-destination pairs/users and propose a pricing framework that enables the relay to set prices to maximize either its revenue or any desirable system utility. Specifically, depending on the quality of the received signals, the relay sets prices and correspondingly charges the users utilizing its resources for their transmissions. The price is determined in such a way that the relay's revenue or system utility is maximized. Given the specified price, the users competitively employ the relay node to forward their signals. We model each user as a rational player, which aims at maximizing its own net utility through power allocation, and analyze the competition among the users within the framework of non-cooperative game theory. It is shown that, in the game played by the users, there always exists a unique pure Nash equilibrium point that can be achieved through distributed iterations. Next, subject to the availability of complete information about the users at the relay, we propose a low-complexity uniform pricing algorithm and an optimal differentiated pricing algorithm, in which the relay either charges the users at a sub-optimal uniform price or it charges different prices per user. We also show that, by applying the differentiated pricing algorithm that enforces the users to transmit at certain power levels, any system utility can be maximized. Extensive simulations are conducted to quantify the performance of the proposed methods.

## 2.1 Introduction

For many wireless networks, the transmission between two distant users may have to be accomplished with the help of an intermediate node, i.e., relay, due to transmit power or other constraints [1]. In the presence of a relay node, distributed spatial diversity, or cooperative diversity, can be created without physically packing multiple antennas into small-size nodes as long as certain signal combining techniques are applied at the destination [2][3].

The traditional network resource allocation largely relies on system-wide centralized management, which requires all the users to cooperatively follow the resource sharing mechanism and incurs a heavy spectral loss due to the signalling overhead associated with the information exchange and coordination. Nevertheless, in the absence of a central controller, rational or selfish users have incentives to optimize their own performances independently, without considering the social welfare and thus, the existing centralized mechanisms are no longer applicable in such settings. An alternative solution is to model a network of selfish users using non-cooperative game theory [34]. Furthermore, it has been demonstrated in the literature that appropriate pricing techniques can be deployed among multiple selfish users to implement various resource allocation policies, including, but not limited to, revenue maximization [9], social-welfare improvement [18], user fairness guarantee [20] and system-wide optimization [19]. Interested readers are referred to [8] for a survey on game-theoretic resource allocation and pricing mechanisms. In wireless relay networks, without a proper compensation framework, relays have no incentives to forward the signals of various users to the corresponding destinations, since this is done at the expense of their own energy consumption. Hence, pricing becomes a useful and efficient mechanism that reimburses relays for using their resources by making payments,[1] thereby providing the

---

[1]The payments can be tokens, virtual money, etc., which can be used in the future by the relay to purchase resources from the other nodes in the network.

relays with incentives to forward the other users' signals [11]–[13].

In this chapter, we focus on a wireless relay network, in which there exists one relay node and multiple source-destination pairs/users.[2] We propose a pricing mechanism that gives the relay incentives to forward the users' signals to the destinations. In particular, the price is determined by the relay such that its revenue[3] or system utility is maximized. Given the specified price, the users competitively utilize the relay node to forward their signals and make appropriate payments to the relay based on the receive signal to interference plus noise ratio (SINR). We model each user as a selfish player, which aims at maximizing its own net utility by adjusting its transmit power, and analyze the emerging competition among the users using non-cooperative game theory. Specifically, given the knowledge of its local channel state information (CSI), each user maximizes its utility by optimally choosing its power level in response to the power allocation strategies of the other users. This process iterates until convergence. We show that, in the non-cooperative game played by the users, there always exists a unique pure Nash equilibrium point (NEP) that can be achieved through the distributed iterative power allocation process. Next, we assume that the relay has only incomplete information about the users (i.e., the number of users and the sum SINR when all the users transmit with their maximum powers) and propose a low-complexity uniform pricing algorithm based on which the relay charges the users at a sub-optimal uniform price. Subsequently, we extend the uniform pricing algorithm to differentiated pricing by assuming that the relay has complete information about the users (i.e., channel coefficients, power constraints, etc.). Furthermore, we show that, by utilizing the differentiated pricing algorithm, any system utility can be maximized even though

---

[2]Throughout this chapter, we interchangeably use the term "user" to represent the source-destination pair.

[3]The dedicated relay incurs a fixed cost, e.g., power consumption, associated with forwarding the users' signals, and moreover, the relay's resource in the current time slot cannot be reserved for further use [17]. For instance, the cost of deploying the relay station and power expenditure is paid in advance by the infrastructure manager. Therefore, as described in the transmission protocol, the relay will forward the users' signals and revenue maximization is virtually equivalent to profit maximization [25][27][33].

the users behave selfishly. Finally, extensive simulations are conducted to verify the performance of the proposed methods.

The main contributions of this chapter are threefold: (i) we focus on a relay network with multiple users that are modeled as selfish players competing against each other for the scarce network resource, i.e., the relay, and study the NEP of the non-cooperative game; (ii) depending on how much information the relay has about the users, we propose two pricing algorithms, i.e., uniform pricing with incomplete information and differentiated pricing with complete information; (iii) the proposed differentiated pricing algorithm enforces the users to transmit at desired power levels at the NEP and hence, can be applied to optimize any system utility, which includes the relay's revenue as a special case.

The rest of this chapter is organized as follows. Chapter 2.2 provides the literature review and Chapter 2.3 describes the system model and problem formulation. In Chapter 2.4, a distributed power allocation algorithm along with two pricing algorithms are developed for the considered relay network. Simulation results are shown in Chapter 2.5. Finally, concluding remarks are offered in Chapter 2.6.

## 2.2   Related Works

Power allocation, both with and without pricing, has been extensively studied in wireless networks. Next, we present a brief overview of the related works and describe the relationship to our proposed mechanism.

It is worth noting that pricing mechanisms, which originate from the competitive market theory [33], have been widely applied in the context of cognitive systems [9][10] and relay networks [11]–[13]. For instance, given a wide-band uplink cognitive system, [9] proposes a differentiated pricing algorithm that charges different secondary users at different prices to maximize the revenue of the service provider.

To utilize the benefits of distributed spatial diversity and guarantee incentive compatibility in wireless cooperative networks, [10] adopts the hierarchical Stackelberg game-theoretic framework where the primary user, as the leader, selects some secondary users as the cooperative relay nodes and, in return, grants the spectrum usage to the participating secondary users for their own data transmissions. As followers, the secondary users decide the payment made to the primary user to gain the channel access time and maximize their own utilities. Considering a cooperative network with multiple self-interested relays, the authors in [11] cast the problem of distributed power control and relay selection into the Stackelberg formulation. In particular, the relays are regarded as leaders that selfishly set the prices such that they can maximize the revenue. The payment made by the user serves as a reimbursement that gives the relay an incentive to participate in the cooperation. Similar compensation frameworks enabling the relay to forward the users' signals are proposed in the literature, e.g., [12][13].

Following a joint user-centric and network-centric optimization approach, the authors in [14] propose a distributed power control and revenue optimization framework in conventional cellular networks. Specifically, the network controller, e.g., base station, charges each user in accordance with its throughput while the users transmit over an interference channel and maximize the energy efficiency. The same approach is later applied in the multi-cell scenario [15]. In [16], an auction-based spectrum sharing protocol is proposed such that each user submits an optimal bid to the network manager to maximize the utility minus the payment. Two payment rules, i.e., SINR and power, are considered and it is shown that, with logarithmic utilities, the power auction outperforms the SINR auction in terms of the revenue from the network perspective. The auction framework is also extended in [17] to a cooperative network setting wherein the relay and the users are modeled as the auctioneer and bidders, respectively. Focusing on the classic Gaussian interference channel, [24] introduces

the notion of "taxation" which captures the effect of one user's power allocation on the others', and presents a modified iterative water-filling algorithm that maximizes the sum utility. For a cellular network, the authors in [26] proposed a differentiated pricing mechanism such that any near-optimal system utility can be achieved.

In contrast with the existing literature, we shift our attention to a relay network with multiple selfish users and propose a pricing mechanism that can maximize the system utility and provide the relay with incentives to forward the signals of the users. New challenges emerging in such relay networks include: (i) how to design a proper pricing mechanism and how to set the price; (ii) given the price, how to model and analyze the competition among the selfish users; (iii) in view of the users' selfishness, how to maximize the system utility. In this chapter, the users competitively adjust their transmit powers and utilize the relay node to accomplish their own transmissions and, as the service provider, the relay charges all the users according to either uniform or differentiated pricing algorithm to maximize the revenue. Furthermore, the differentiated pricing can be applied to optimize any system utility.

## 2.3   System Model

Consider a wireless relay network consisting of one relay node and multiple source-destination pairs.[4] The sources and destinations are indexed by $\mathcal{S}_i$ and $\mathcal{D}_i$, respectively, for $i = 1, 2 \cdots Q$, and the relay node is represented by $\mathcal{R}$. Similar system models with one relay node and multiple users have been considered in the literature for different purposes as well [4]–[7], wherein [4]–[6] focus on orthogonal transmissions without interferences and [7] studies interference management in a two-way relay channel [1].

---

[4]As in [4][5], the analysis throughout this chapter can be applied to a network with more than one relays, provided that the network can be classified into multiple clusters, each of which consists of one relay and multiple users, and different clusters are transmitting over different channels.

### 2.3.1 Network Model

We assume that the channels are flat (or frequency non-selective) fading. When the channels are frequency-selective fading and divided into multiple subchannels (e.g., OFDM), the proposed algorithm in this chapter can still be applied on a per-subchannel basis if each user has an individual maximum power constraint for each subchannel. Nevertheless, if each user has a total power constraint across all the subchannels, it is intrinsically difficult to generalize the proposed algorithm. In order to keep the analysis tractable, we note that it is a common practice to focus on a flat fading (or frequency non-selective fading) channel model when studying pricing-related algorithms ([11][14][17][26]). The channel coefficients for the $\mathcal{S}_i - \mathcal{R}$ and the $\mathcal{R} - \mathcal{D}_i$ channels are denoted by $g_i$ and $h_i$, respectively, for $i = 1, 2 \cdots Q$. The transmit powers of $\mathcal{S}_i$ and $\mathcal{R}$ are $p_i$ and $p_{\mathcal{R}}$, respectively. The local CSI, i.e., $g_i$ and $h_i$, is only obtained by user $i$, and neither $g_j$ nor $h_j$ is known to user $i$, if $j \neq i$, due to the distributed nature of the considered communication problem. Furthermore, we assume that the zero-mean complex additive white Gaussian noise (AWGN) at each node to has a variance[5] of $N_0$. Due to the half-duplex constraint, we consider orthogonal relaying transmissions, e.g., the source nodes and the relay node transmit in two non-overlapping time slots. The direct link between $\mathcal{S}_i$ and $\mathcal{D}_i$ is neglected due to, for instance, the shadowing effects [1]. To forward the data from the source to the destination, we adopt the classical amplify-and-forward strategy [3] as the relaying operation, which has been shown to be an appealing technique due to its low cost and easy implementation as compared to the decode-and-forward protocol [30]. Hence, the signals received at $\mathcal{R}$ and $\mathcal{D}_i$ can be written, respectively, as

$$y_{\mathcal{R}} = \sum_{j=1}^{Q} g_j \sqrt{p_j} x_j + n_{\mathcal{R}} \text{ and } y_i = \alpha h_i y_{\mathcal{R}} + n_i, \qquad (2.1)$$

---

[5]This assumption is imposed only for the convenience of notation, as in [24], and can be relaxed without affecting the analysis in this chapter.

where $x_i$ is the unit-variance transmit signal from $\mathcal{S}_i$ to $\mathcal{D}_i$, $\alpha$ is the amplification factor of $\mathcal{R}$, $n_{\mathcal{R}}$ and $n_i$ are the statistically-independent AWGN terms at $\mathcal{R}$ and $\mathcal{D}_i$, respectively. The amplification factor $\alpha = \sqrt{\frac{p_{\mathcal{R}}}{\sum_{j=1}^{Q}|g_j|^2 p_j + N_0}}$, which is public information available to all the users, is chosen to satisfy the fixed power constraint at the relay. Assuming that $\mathcal{D}_i$ is only interested in the signal $x_i$ and treats the multi-user interference as noise [23][29], we can then express the receive SINR at $\mathcal{D}_i$ as

$$\gamma_i = \frac{|g_i|^2|h_i|^2 p_{\mathcal{R}} p_i}{|g_i|^2 N_0 p_i + (|h_i|^2 p_{\mathcal{R}} + N_0) \cdot \left(\sum_{j=1, j \neq i}^{Q} |g_j|^2 p_j + N_0\right)}. \tag{2.2}$$

Recall that in an amplify-and-forward relay network with only one source node, only AWGN noise is amplified and forwarded by the relay to the destination node in addition to the desired signal component, and thus, the received signal to noise ratio is expressed as $\frac{|g_i|^2|h_i|^2 p_{\mathcal{R}} p_i}{|g_i|^2 N_0 p_i + (|h_i|^2 p_{\mathcal{R}} + N_0) N_0}$. When there are multiple source nodes transmitting simultaneously to the relay, both AWGN noise and multi-user interference (i.e., $\sum_{j=1, j \neq i}^{Q} |g_j|^2 p_j$) are amplified and forwarded to the destination nodes and hence, the resulting SINR expression becomes (2.2). The achievable rate of user $i$ is therefore given by

$$R_i(p_i; p_{-i}) = \frac{1}{2}\log\left(1 + \gamma_i\right), \tag{2.3}$$

where the scaling factor $1/2$ is due to the fact that $\mathcal{S}_i$ transmits $x_i$ only for half of the frame, $\gamma_i$ is given in (2.2), and $p_{-i} = (p_1 \cdots p_{i-1}, p_{i+1} \cdots p_Q)$ is the vector of power allocation strategies of all the users except for user $i$.

Before proceeding to the problem formulation, we briefly discuss how transmissions using the relay node considered in this chapter significantly differ from conventional single-hop transmissions [23][24], despite the absence of direct channels. First, the signals are transmitted through a cascaded channel, i.e., multi-access channel followed by broadcast fading channel. Second, the signal forwarded by the relay

node is not "clean", whereas the source transmits noiseless signals to the destination in single-hop Gaussian interference channels, i.e., the relay amplifies the Gaussian noise, in addition to the desired signal, which can be seen from the signal model in (2.1). Hence, the analysis in this chapter can be regarded as a generalization of the existing results on one-hop interference channels. As a special case, if the relay-destination channel is sufficiently good (i.e., $|h_i|^2 \to \infty$), the dual-hop relay channel reduces to the conventional multi-access interference channel and the receive SINR of user $i$ becomes $\gamma_i = \frac{|g_i|^2 p_i}{\sum_{j=1, j \neq i}^{Q} |g_j|^2 p_j + N_0}$, which can also be obtained by taking the limit of (2.2) with respect to $|h_i|^2 \to \infty$.

### 2.3.2 Problem Formulation

There are various payment rules in communications networks. For instance, each individual user may be charged in proportion to the relay's transmit power [11], its throughput [14], receive SINR [16][17], allocated rate [25], and its own transmit power [26]. In the problem considered in this chapter, it is clear from (2.2) that the receive SINR is partially determined by the relay's power. Furthermore, it is the SINR that measures the quality of the received signal and thus influences the utility of each user. Hence, it is reasonable to assume that the payment made to the relay is a function of the receive SINR. In particular, we assume in this chapter that the payment that user $i$ needs to make to the relay is defined as $\pi_i \gamma_i$, where $\pi_i$ is the price for user $i$ set by the relay. As will be shown in this chapter, this payment rule allows the relay to set optimal differentiated prices such that any system utility function can be maximized. Moreover, the considered payment rule charges each user in proportion to its receive SINR, and has been applied in [16] and [17] for different purposes (e.g., to achieve different tradeoffs between fairness and efficiency in a multi-user relay network [17]). Other similar payment rules can be found in [9][14]. In general, the utility function is increasing and concave in the receive SINR [16][17]. As a particular example and

to gain more insights on the pricing algorithms, we adopt in the sequel the achievable rate[6] $R_i(p_i; p_{-i})$ as the utility function of user $i$. Given the utility function and payment rule, the payoff, or net utility function, of user $i$ can therefore be expressed as the following surplus [9][11][17]

$$u_i(p_i; p_{-i}) = \frac{1}{2}\log\left(1 + \gamma_i\right) - \pi_i\gamma_i, \tag{2.4}$$

where the first term $\frac{1}{2}\log\left(1 + \gamma_i\right)$ is the achievable rate of user $i$, and $\pi_i\gamma_i$ is the payment made to the relay. From the perspective of the relay, in order to maximize the revenue collected from the users, the relay needs to set an optimal price vector $\mathbf{\Pi}^* = \left\{\pi_1^*, \pi_2^* \cdots \pi_Q^*\right\}$ such that

$$\mathbf{\Pi}^* = \arg\max_{\mathbf{\Pi} \succeq \mathbf{0}} \left(\sum_{i=1}^{Q} \pi_i\gamma_i(p_i; p_{-i})\right). \tag{2.5}$$

Note that the relay's price and the users' power allocation are coupled in a complex way, and we shall address the coupling in Chapter 2.4. In particular, relay's price influences the users' power allocation which, in turn, affects the relay's revenue. While we first use the revenue as a particular utility for the relay, we note that the proposed pricing mechanism can also be applied to maximize any system utility, making the proposed pricing framework a suitable option for managing wireless relay networks with selfish users.

---

[6]Note that the logarithm-based function or achievable rate is a widely-used utility definition (see, e.g., [9][11][17][24][25]) and the analysis herein can be applied, after modifications, to other concave utility functions as well. In particular, the existence of pure NEP and convergence of the iterative power control algorithm are not affected if we replace the achievable rate with a general concave utility function in (2.4). Furthermore, it is easy to incorporate a *weight* into the utility function, i.e., user $i$ has a utility of $w_i R_i(p_i; p_{-i})$ where $w_i > 0$ is a factor that converts the achievable rate into currency [9] or approximates the reception quality in the case of video delivery applications [28].

## 2.4 User-Centric Optimization and Pricing

In this part, we cast the user-level problem of distributed power allocation into the framework of non-cooperative game theory. Adopting revenue as the relay's utility, we propose two pricing algorithms, i.e., uniform pricing and differentiated pricing. Then, we show that the differentiated pricing algorithm can maximize any system utility by enforcing the users to transmit at desired power levels.

### 2.4.1 Distributed Power Allocation

Non-cooperative game theory is an effective tool to capture the selfish behaviors of self-interested players [34]. Given the price set by the relay, we can mathematically characterize the competition among the self-interested users as a non-cooperative game

$$\mathcal{G}_{\text{user}} = \{\Omega, \{\mathcal{P}_i\}_{i \in \Omega}, \{u_i(p_i; p_{-i})\}_{i \in \Omega}\} \tag{2.6}$$

where $\Omega \triangleq \{1, 2 \cdots Q\}$ is the set of active users (i.e., $\mathcal{S}_i - \mathcal{D}_i$ pair), $\mathcal{P}_i$ is the set of admissible power allocation strategies of user $i$ defined as $\{p_i : 0 \leq p_i \leq p_i^{\max}\}$ and $u_i(p_i; p_{-i})$ is the payoff of user $i$ given in (2.4). The optimal power of user $i$ in response to the power levels of all the other users is referred to as the *best response* function denoted by $\mathcal{B}_i(p_{-i})$. In the non-cooperative game played by the users, the NEP is achieved when user $i$, given $p_{-i}$, cannot increase its net utility $u_i(p_i; p_{-i})$ by unilaterally changing its own power $p_i$, for all $i \in \Omega$. Mathematically, the NEP, denoted by $\mathbf{p}^* = (p_1^*, p_2^* \cdots p_Q^*)$, of the user-level game $\mathcal{G}_{\text{user}}$ in (2.6) is formally defined as follows [34]

$$u_i(p_i^*; p_{-i}^*) \geq u_i(p_i; p_{-i}^*), \quad \forall p_i \in \mathcal{P}_i, \ \forall i \in \Omega. \tag{2.7}$$

It is known that, in a one-shot[7] non-cooperative game, pure NEP is a critical operating point at which the outcome of the game becomes stabilized [34], and thus, it is of great interest to study the existence of NEP in such a game. Moreover, whether and how the non-cooperative game can eventually arrive at the NEP is another question we have yet to answer. To this end, we first explicitly express the best response function of user $i$, i.e., $\mathcal{B}_i(p_{-i})$, which specifies the transmit power user $i$ should use in response to the other users' power strategies and the price set by the relay. Specifically, depending on the price $\pi_i$ set by the relay, the unique $\mathcal{B}_i(p_{-i})$ can be derived and expressed in a compact form as

$$\mathcal{B}_i(p_{-i}) = \left[ \frac{\delta_i(\pi_i)\,(|h_i|^2 p_{\mathcal{R}} + N_0)\left(\sum_{j=1, j\neq i}^{Q} |g_j|^2 p_j + N_0\right)}{|g_i|^2 \cdot \left[|h_i|^2 p_{\mathcal{R}} - N_0 \cdot \delta_i(\pi_i)\right]} \right]_0^{p_i^{\max}} \tag{2.8}$$

where $\left[\ \cdot\ \right]_a^b = \max\{\min\{\ \cdot\ , b\}, a\}$ and $\delta_i(\pi_i)$ is a non-negative and continuously non-increasing function of $\pi_i$ defined as

$$\delta_i(\pi_i) = \begin{cases} 0, & \text{if } \frac{1}{2} < \pi_i, \\ \frac{1}{2\pi_i} - 1, & \text{if } (1 + \gamma_i(p_i^{\max}; \mathbf{0}))^{-1} < 2\pi_i \leq 1, \\ \gamma_i(p_i^{\max}; \mathbf{0}), & 0 \leq 2\pi_i \leq (1 + \gamma_i(p_i^{\max}; \mathbf{0}))^{-1}, \end{cases} \tag{2.9}$$

in which $\gamma_i(p_i^{\max}; \mathbf{0})$ is obtained by plugging $(p_i; p_{-i}) = (p_i^{\max}; \mathbf{0})$ into (2.2). Denote $\mathbf{p} = (p_1, p_2 \cdots p_Q)$ and $\mathcal{B}(\mathbf{p}) = (\mathcal{B}_1(p_{-1}), \mathcal{B}_2(p_{-2}) \cdots \mathcal{B}_Q(p_{-Q}))$, respectively. Then, in order to facilitate the analysis and development of the distributed algorithm, we further simplify (2.8) and express it in a vector form as

$$\mathcal{B}(\mathbf{p}) = [\mathbf{T}\,\mathbf{p} + \mathbf{t_0} N_0]_{\mathbf{0}}^{\mathbf{p}^{\max}}, \tag{2.10}$$

---

[7]As will be shown later, the pure NEP is reached through an iterative power allocation process. Nevertheless, the user-level game in this chapter is still one-shot in the sense that, unlike in a repeated game [34], the players or users do not take into account the history or future utility when making the current decisions. Thus, pure NEP is an appropriate concept characterizing the steady outcome of the game.

where $[\mathbf{a}]_{\mathbf{0}}^{\mathbf{p}^{\max}} = \left([a_1]_0^{p_1^{\max}}, [a_2]_0^{p_2^{\max}} \cdots [a_Q]_0^{p_Q^{\max}}\right)$, $\mathbf{t_0}_i = \frac{\delta(\pi_i)\cdot\left(|h_i|^2 p_\mathcal{R} + N_0\right)}{|g_i|^2\cdot[|h_i|^2 p_\mathcal{R} - N_0\cdot\delta(\pi_i)]}$ and

$$\mathbf{T}_{ij} = \begin{cases} \frac{\delta(\pi_i)\cdot\left(|h_i|^2 p_\mathcal{R} + N_0\right)}{|g_i|^2\cdot[|h_i|^2 p_\mathcal{R} - N_0\cdot\delta(\pi_i)]} \cdot |g_j|^2, & \text{if } i \neq j, \\ 0, & \text{if } i = j, \end{cases} \tag{2.11}$$

for $i, j = 1, 2 \cdots Q$. Next, we present an iterative distributed algorithm (i.e., Algorithm I), in which each user chooses, at each iteration, its best response to the power strategies of the others.

### Algorithm I: Iterative Distributed Power Allocation

. . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .

**Input:**    $\pi_i, g_i, h_i$ **for user** $i$, $i = 1, 2 \cdots Q$
**Step 1:**   $n = 0$; **choose any feasible** $\mathbf{p}^0 = \left(p_1^0, p_2^0 \cdots p_Q^0\right)$
**Step 2:**   $\mathbf{p}^{(n+1)} = \mathcal{B}(\mathbf{p}^n)$
**Step 3:**   $n = n + 1$; **go to Step 2 until convergence**

. . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .

In Step 1 of Algorithm I, each user $i$ can arbitrarily choose its initial power $p_i^0$ from its feasible power set $\{p_i : 0 \leq p_i \leq p_i^{\max}\}$ in a distributed manner and then, the initial power vector $\mathbf{p}^0 = (p_1^0, p_2^0 \cdots p_Q^0)$ is also feasible. To complete the algorithm description, we give Theorem 1 regarding the existence of NEP in the user game and the convergence of the proposed algorithm.

**Theorem 1.** *Given any non-negative price vector* $\mathbf{\Pi} \succeq \mathbf{0}$ *set by the relay, there exists a unique pure NEP of the user game* $\mathcal{G}_{\text{user}}$. *Moreover, starting from any initial point* $\mathbf{p}^0 \in \mathcal{P} \triangleq \mathcal{P}_1 \times \mathcal{P}_2 \cdots \times \mathcal{P}_Q$, *the iteration specified by* $\mathbf{p}^{(n+1)} = \mathcal{B}(\mathbf{p}^n)$ *always converges to the unique NEP of the user game as* $n \to \infty$.

*Proof.* The proof is mainly based on the *standard* interference function that was first proposed for distributed power control in [21]. Any function $f(\mathbf{x})$ satisfying the following three properties, for all $\mathbf{x} \geq \mathbf{0}$, is called *standard*:

1. Positivity: $f(\mathbf{x}) > 0$;

20

2. Monotonicity: if $\mathbf{x} \geq \tilde{\mathbf{x}}$, then $f(\mathbf{x}) \geq f(\tilde{\mathbf{x}})$;

3. Scalability: for all $\beta > 1$, $\beta f(\mathbf{x}) > f(\beta \mathbf{x})$.

To prove the existence of a unique NEP and the convergence of Algorithm I, we consider the following two cases depending on the value of $\mathbf{\Pi}$ which plays a critical role in the best response vector $\mathcal{B}(\mathbf{p})$.

*Case 1:* $\mathbf{0} \preceq \mathbf{\Pi} \prec \frac{1}{2}$.

We have in this case $0 < \mathbf{t_{0i}} < +\infty$ and $0 < \mathbf{T}_{ij} < +\infty$, for $i \neq j$ and $i, j = 1, 2 \cdots Q$. It is trivial to show that, without considering the maximum power constraint, the function of $\mathbf{p}$ in the update (2.10), i.e., $\mathbf{T} \mathbf{p} + \mathbf{t_0} N_0$, is standard. Then, following the proof of Theorem 7 in [21], we can easily prove that the update rule defined in (2.10) with the maximum power constraint is also standard. Hence, by applying Corollary 1 in [21], we establish the existence of a unique fixed point in the proposed iterative power allocation process (i.e., NEP of the user game) and the convergence of Algorithm I to this unique NEP, given $\mathbf{0} \preceq \mathbf{\Pi} \prec \frac{1}{2}$.

*Case 2:* $\frac{1}{2} \leq \pi_i$, for some $i \in \Omega$.

In this case, the iteration $p_i^{n+1} = \mathcal{B}_i(\mathbf{p}^n)$ is always zero, if $\frac{1}{2} \leq \pi_i$. Thus, users that are charged a price greater than or equal to $\frac{1}{2}$ can be excluded from the network. The remaining users are all charged with a price less than $\frac{1}{2}$ and hence, they form a new *virtual* network that satisfies *Case 1*. Hence, as we have shown in *Case 1*, the game played by the users in the virtual network admits a unique NEP that can be reached by applying Algorithm I. Note that adding the users that are charged a price greater than or equal to $\frac{1}{2}$ into the virtual network has no effect for the virtual network, since the added users always transmit at zero powers. Therefore, the game has a unique NEP and the proposed distributed power allocation algorithm converges to this unique NEP regardless of the initial point, even though $\mathcal{B}(\mathbf{p})$ is not a standard interference function as it violates the properties of *positivity* and *scalability*.

To sum up, we have proved Theorem 1 by considering the above two cases. The existence of a NEP can also be proved by showing that the net utility function of each user is quasi-concave in this user's power and continuous in the power of all the users, and that the feasible power set is compact and closed. The details are omitted for brevity. It should also be noted that, in general, the existence of a (even unique) fixed point of an iterative process does not necessarily imply the convergence of this iterative process (see [35] for an example). The existence of a fixed point and convergence are two separate properties of an iterative process. In the problem considered in this chapter, however, both the existence of a fixed point (NEP) and the convergence of the iterative process can be established, since the best response function is standard and there exists a maximum power constraint [21]. ∎

Before concluding this part, we note that the distributed nature of the algorithm stems from the fact that the information required to compute $\mathcal{B}_i(p_{-i})$ at user $i$ can be locally observed without exchanging CSI among different users. Specifically, as shown in (2.8), the information needed by user $i$ includes the local CSI (i.e., $g_i$ and $h_i$), the relay's transmit power $p_{\mathcal{R}}$, the price $\pi_i$ set by the relay and the multi-user interference plus noise $\sum_{j=1, j\neq i}^{Q} |g_j|^2 p_j + N_0$. In particular, user $i$ can obtain the local CSI through channel estimation and feedback.[8] The relay's transmit power $p_{\mathcal{R}}$ and the price $\pi_i$ are transmitted via control channels to user $i$ prior to the users' transmissions. Regarding the multi-user interference, the relay node can broadcast to all the users its amplification factor $\alpha$ such that user $i$, for $i \in \Omega$, acquires the value of $\sum_{j=1, j\neq i}^{Q} |g_j|^2 p_j + N_0$ by computing $\frac{p_{\mathcal{R}}}{\alpha^2} - |g_i|^2 p_i = \sum_{j=1, j\neq i}^{Q} |g_j|^2 p_j + N_0$. It can therefore be seen that the proposed algorithm can be applied in a distributed manner and that it needs to be re-executed when the price set by the relay is updated or the network condition changes. Finally, note that the proposed algorithm is applicable to scenarios in which the en-

---

[8]At the beginning of a frame, a known pilot symbol is sent by a transmitter node to allow its receiver node to estimate the channel gain and then feed it back to the transmitter node. Other schemes are also available to let the users obtain their local CSI (see [32] and references therein for details).

vironment does not change frequently (e.g., the channel gains vary slowly when the nodes in the network move sufficiently slowly or remain in fixed positions). On the other hand, if the channels suffer from fast fading (e.g., due to high mobility), the proposed algorithm no longer works. The same limitation exists in (almost) all the existing work (see, for instance, [11][16][17][24]).

### 2.4.2 Uniform Pricing With Incomplete Information

In many wireless networks with limited information exchange among different nodes, the relay has only incomplete information about the users (e.g., the maximum power constraints of the users are private and thus unknown to the relay). Under such constraints, we propose a uniform pricing algorithm, i.e., the relay sets and broadcasts to the users a uniform price $\pi_1 = \pi_2 \cdots = \pi_Q = \pi$.

As we have stated in Theorem 1, the user-level game always consists of a unique NEP given any price vector set by the relay. Hence, the relay aims at maximizing its revenue by setting an appropriate price when the game reaches the NEP, i.e., the user game becomes stabilized. Nevertheless, since the private information of the users, e.g., power strategy space, is unknown to the relay, it cannot analytically compute the NEP of the user-level game $\mathcal{G}_{\text{user}}$ or directly set an optimal uniform price such that $\pi^* = \arg\max_{\pi \geq 0} \left( \pi \sum_i^Q \gamma_i(p_i^*; p_{-i}^*) \right)$. As a consequence, an iterative process that adjusts the price is needed to identify the optimal uniform price. A naive solution is to perform brute-force exhaustive search. Specifically, the relay divides the range of feasible prices into many sufficiently small intervals, and for each small interval, the relay selects a uniform price that falls into the interval and computes the revenue when the user-level game reaches the unique NEP. Finally, the relay chooses the price that generates the maximum revenue among all the candidate prices. Unfortunately, the average total number of iterations required by this method to obtain the optimal uniform price is $m\bar{N}$, where $m$ is the number of candidate prices and is typically a

large value, and $\bar{N}$ is the average number of iterations needed by the distributed power allocation algorithm to converge.

Given that it is computationally prohibitive and mathematically involved to find the optimal uniform price through the exhaustive search, we alternatively propose a low-complexity algorithm that can yield a close-to-optimal uniform price. Before stating the algorithm, we first define the lower and upper bounds on the optimal uniform price, i.e., $\pi_a = \frac{1}{2} \min_{i \in \Omega} \{1 + \gamma_i (\mathbf{p}^{\max})\}^{-1}$ and $\pi_b = \frac{1}{2} \max_{i \in \Omega} \{1 + \gamma_i (\mathbf{p}^{\max})\}^{-1}$, respectively, and summarize some instrumental properties of the revenue function[9] $\rho(\pi) = \pi \sum_{i=1}^{Q} \gamma_i(\pi)$ in the following theorem.

**Theorem 2.** *The revenue function has the following properties:*

*1. $\rho(\pi) \geq 0$;*

*2. $\rho(\pi) = 0$ if $\pi = 0$ or $\pi \geq \frac{1}{2}$;*

*3. $\rho(\pi) < \infty$ if the number of users, $Q$, is finite.*

*4. $\rho(\pi) = \pi \sum_{i=1}^{Q} \gamma_i(\mathbf{p}^{\max})$ when $0 \leq \pi \leq \pi_a$;*

*5. There exists a certain value of price $\hat{\pi}$ satisfying*

$$\begin{cases} \hat{\pi} < \pi_b, & \exists i, j \in \Omega \ s.t. \ \gamma_i (\mathbf{p}^{\max}) \neq \gamma_j (\mathbf{p}^{\max}) \\ \hat{\pi} = \pi_b, & \forall i, j \in \Omega \ s.t. \ \gamma_i (\mathbf{p}^{\max}) = \gamma_j (\mathbf{p}^{\max}) \end{cases}, \tag{2.12}$$

*such that $\rho(\pi) = Q \cdot \left(\frac{1}{2} - \pi\right)$ if $\hat{\pi} \leq \pi \leq \frac{1}{2}$;*

*Proof.* The proof is given in the order of the properties listed in Theorem 2.

Property 1–3 directly follows the best response function in (2.8).

Property 4: Given $p_{-i} = \left(p_1^{\max} \cdots p_{i-1}^{\max}, p_{i+1}^{\max} \cdots p_Q^{\max}\right)$, it can be derived from the best response function that $\mathcal{B}_i(p_{-i}) = p_i^{\max}$, if $0 \leq \pi \leq \frac{1}{2} \{1 + \gamma_i(\mathbf{p}^{\max})\}^{-1}$. Hence, $\mathbf{p}^{\max} = \left(p_1^{\max}, p_2^{\max} \cdots p_Q^{\max}\right)$ satisfies $\mathbf{p}^{\max} = \mathcal{B}(\mathbf{p}^{\max})$, i.e., $\mathbf{p}^{\max}$ is the NEP,

---

[9]The SINR is an explicit function of the uniform price $\pi$ which affects the net utility and the power allocation of users.

when $0 \le \pi \le \frac{1}{2} \min_{i=1,2\cdots Q} \left\{ 1 + \gamma_i \left( \mathbf{p}^{\max} \right) \right\}^{-1}$. In this case, by Theorem 2, the distributed power allocation algorithm globally converges to the unique point $\mathbf{p}^{\max}$. The intuitive interpretation is that, when the price is sufficiently low, every user can afford the payment charged by the relay and thus will transmit at a high power. When all the users transmit at their maximum powers, the receive SINR $\gamma_i(\pi)$ is a positive constant, denoted by $\tilde{\gamma}_i = \gamma_i(\mathbf{p}^{\max})$, for $i = 1, 2 \cdots Q$, irrespective the value of $\pi$. Therefore, the revenue $\rho(\pi) = \pi \sum_j^Q \tilde{\gamma}_j$ is a strictly increasing function of $\pi$ when $0 \le \pi \le \frac{1}{2} \min_{i=1,2\cdots Q} \left\{ 1 + \gamma_i \left( \mathbf{p}^{\max} \right) \right\}^{-1}$.

Property 5: We first introduce the following lemma before proving the existence of $\hat{\pi}$.

**Lemma 1.** If $\pi > \frac{1}{2} \left\{ 1 + \gamma_i(\mathbf{p}^{\max}) \right\}^{-1}$, then the maximum transmit power constraint of user $i$ is not activated at the NEP of the game $\mathcal{G}_{\text{user}}$, i.e., $0 \le p_i^* < p_i^{\max}$, for any $i = 1, 2 \cdots Q$.

*Proof.* By taking the first-order derivative of $\mathcal{B}_i(p_{-i})$ in (2.8) with respect to $\pi$, it can be easily shown that $\mathcal{B}_i(p_{-i})$ is a strictly decreasing function of $\pi$ when $\frac{1}{2} \left\{ 1 + \gamma_i(p_i^{\max}; p_{-i}) \right\}^{-1} \le \pi \le \frac{1}{2}$. In particular, $\pi = \frac{1}{2} \left\{ 1 + \gamma_i(p_i^{\max}; p_{-i}) \right\}^{-1}$ results in $\mathcal{B}_i(p_{-i}) = p_i^{\max}$. Therefore, the maximum power constraint of user $i$ is not activated, i.e., $0 \le \mathcal{B}_i(p_{-i}) < p_i^{\max}$, if $\frac{1}{2} \left\{ 1 + \gamma_i(\mathbf{p}^{\max}) \right\}^{-1} < \pi \le \frac{1}{2}$.

When $\pi > \frac{1}{2}$, $\mathcal{B}_i(p_{-i})$ is always zero. Hence, Lemma 1 is proved. □

Now, we shall prove Property 5 by considering the following two cases.

*Case 1:* $\gamma_i \left( \mathbf{p}^{\max} \right) = \gamma_j \left( \mathbf{p}^{\max} \right)$, for $i, j = 1, 2 \cdots Q$.

In this case, we will show that $\hat{\pi} = \frac{1}{2} \max_{i=1,2\cdots Q} \left\{ 1 + \gamma_i \left( \mathbf{p}^{\max} \right) \right\}^{-1}$. When $\hat{\pi} \le \pi \le \frac{1}{2}$, $\mathcal{B}_i(p_{-i})$ satisfies $\frac{\partial u_i(p_i; p_{-i})}{\partial p_i} \big|_{p_i = \mathcal{B}_i(p_{-i})} = 0$, for $i = 1, 2 \cdots Q$. Then, it can be derived that

$$\gamma_i \left[ \mathcal{B}_i(p_{-i}); p_{-i} \right] = \frac{1}{2\pi} - 1, \ \forall p_{-i} \in \mathcal{P}_1 \times \cdots \mathcal{P}_{i-1} \times \mathcal{P}_{i+1} \cdots \mathcal{P}_Q. \tag{2.13}$$

Thus, at the NEP of the game $\mathcal{G}_{\text{user}}$, we have $\gamma_i(p_i^*; p_{-i}^*) = \frac{1}{2\pi} - 1$, for $i = 1, 2 \cdots Q$. Therefore, the revenue at the relay, i.e., $\rho(\pi) = \pi \sum_{i=1}^{Q} \gamma_i(p_i^*; p_{-i}^*) = Q \cdot \left(\frac{1}{2} - \pi\right)$, is a strictly decreasing function of $\pi$ when $\hat{\pi} = \frac{1}{2} \max_{i=1,2\cdots Q} \{1 + \gamma_i(\mathbf{p}^{\text{max}})\}^{-1}$.

*Case 2:* "$\gamma_i(\mathbf{p}^{\text{max}}) = \gamma_j(\mathbf{p}^{\text{max}})$, for $i, j = 1, 2 \cdots Q$" does not hold.

Without loss of generality, we assume $\{1 + \gamma_1(\mathbf{p}^{\text{max}})\}^{-1} = \max_{i=1,2\cdots Q} \{1 + \gamma_i(\mathbf{p}^{\text{max}})\}^{-1}$ and $\{1 + \gamma_Q(\mathbf{p}^{\text{max}})\}^{-1} = \min_{i=1,2\cdots Q} \{1 + \gamma_i(\mathbf{p}^{\text{max}})\}^{-1}$. Lemma 1 states that, for any value of the price $\pi > \frac{1}{2}\{1 + \gamma_Q(\mathbf{p}^{\text{max}})\}^{-1}$, the maximum transmit power constraint of user $Q$ is not activated at the NEP, i.e., $p_Q^* < p_Q^{\text{max}}$. Then, following the proof of Lemma 1, it can be also shown that $p_i^* < p_i^{\text{max}}$, for $i = 1, 2 \cdots Q - 1$, if $\pi \geq \frac{1}{2}\{1 + \gamma_1(\mathbf{p}^{\text{max}})\}^{-1}$.

By temporarily relaxing the maximum power constraint, we can express the best response function in (2.10) as

$$\mathbf{p}^{(n+1)} = \mathcal{B}(\mathbf{p}^n) = \mathbf{T}\,\mathbf{p}^n + \mathbf{t_0} N_0, \tag{2.14}$$

where $\mathbf{T}$ is defined in (2.11). It was shown in [22] that, if and only if the spectral radius of $\mathbf{T}$ is less than one, the iteration process specified by (2.14) converges to a unique fixed point, regardless of the initial point, and the fixed point is given by

$$\mathbf{p}^* = \mathbf{T}\,\mathbf{p}^* + \mathbf{t_0} N_0 = (\mathbf{I} - \mathbf{T})^{-1}\,\mathbf{t_0} N_0 = \sum_{i=0}^{\infty} \mathbf{T}^i \mathbf{t_0} N_0. \tag{2.15}$$

As stated above, when $\pi \geq \frac{1}{2}\{1 + \gamma_1(\mathbf{p}^{\text{max}})\}^{-1}$, the transmit power of user $i$ is less than its maximum power constraint, for $i = 1, 2 \cdots Q$, and hence, we have

$$\mathbf{p}^* = [\mathbf{T}\,\mathbf{p}^* + \mathbf{t_0} N_0]_{\mathbf{0}}^{\mathbf{p}^{\text{max}}} = \mathbf{T}\,\mathbf{p}^* + \mathbf{t_0} N_0 = (\mathbf{I} - \mathbf{T})^{-1}\,\mathbf{t_0} N_0 = \sum_{i=0}^{\infty} (\mathbf{T})^i \,\mathbf{t_0} N_0.$$

$$\tag{2.16}$$

It should be noted that, if (2.16) holds, $\gamma_i(p_i^*; p_{-i}^*) = \frac{1}{2\pi} - 1$ can be satisfied at the NEP

of the user game, for $i = 1, 2 \cdots Q$, and as a consequence, the revenue of the relay, i.e., $\rho(\pi) = \pi \sum_{i=1}^{Q} \gamma_i(p_i^*; p_{-i}^*) = Q \cdot \left(\frac{1}{2} - \pi\right)$, is a strictly decreasing function of $\pi$.

As each non-diagonal element of $\mathbf{T}$ is continuously decreasing in $\pi$, it is clear that the transmit power of each user at the NEP, given in terms of the fixed point in (2.15), is also decreasing in $\pi$, if the NEP in the game $\mathcal{G}_{\text{user}}$ without considering the maximum power constraint exists. Thus, the minimum price, denoted by $\hat{\pi}$, which yields a matrix $\mathbf{T}$ with a spectral radius of less than one and satisfies (2.16) must be less than $\frac{1}{2} \left\{1 + \gamma_1\left(\mathbf{p}^{\text{max}}\right)\right\}^{-1}$ and, given the minimum price, only one user reaches its maximum power constraint or multiple (less than $Q$) users reach their corresponding maximum power constraints simultaneously at the NEP. On the other hand, if $\hat{\pi}$ is less than or equal to $\frac{1}{2} \left\{1 + \gamma_Q\left(\mathbf{p}^{\text{max}}\right)\right\}^{-1}$, at least one user will violate the maximum power constraint at the fixed point of the iteration process specified by (2.14) and hence, (2.16) cannot be satisfied. Therefore, $\hat{\pi}$ lies between $\frac{1}{2} \left\{1 + \gamma_Q\left(\mathbf{p}^{\text{max}}\right)\right\}^{-1}$ and $\frac{1}{2} \left\{1 + \gamma_1\left(\mathbf{p}^{\text{max}}\right)\right\}^{-1}$, and when $\hat{\pi} \leq \pi \leq \frac{1}{2}$, the revenue at the relay is a strictly decreasing function of $\pi$.

By considering Case 1 and Case 2 separately, we have proved Property 5.

This proves Theorem 2. ∎

Theorem 2 can be interpreted as follows: <1> The receive SINR is always non-negative and thus, the revenue is also non-negative; <2> The revenue of the relay vanishes when the service of the relay, i.e., packet forwarding, is free or the price is too high; <3> The maximum revenue of the relay is finite as long as the number of users is finite; <4> and <5> The optimal price of the relay lies in a certain interval, i.e., $[\pi_a, \hat{\pi}]$, which depends on the channel conditions and transmit power constraints. Property 4 and 5 significantly reduce the complexity associated with the exhaustive search by eliminating the uniform prices that fall out of the range of the optimal price. They also form the basis of the proposed sub-optimal uniform pricing algorithm. Specifically, the sub-optimal uniform price is obtained by artificially shrinking the interval

(i.e., $[\pi_a, \hat{\pi}]$) to a specific point which is then set as the uniform price. Following these desirable properties of the revenue function, we derive the following corollary.

**Corollary 1.** *There exists an optimal finite uniform price $\pi^*$ such that $\pi_a \leq \pi^* \leq \hat{\pi} \leq \pi_b$, and the corresponding maximum revenue $\rho(\pi)$ is finite and positive. The equalities are activated simultaneously if and only if $\gamma_i\left(\mathbf{p}^{\text{max}}\right) = \gamma_j\left(\mathbf{p}^{\text{max}}\right)$ for all $i, j \in \Omega$.* ∎

Corollary 1 states that the optimal price is upper and lower bounded by $\hat{\pi}$ and $\pi_a$, respectively. As a special case, if and only if $\gamma_i\left(\mathbf{p}^{\text{max}}\right) = \gamma_j\left(\mathbf{p}^{\text{max}}\right)$ for all $i, j \in \Omega$, then we have $\pi_a = \pi^* = \hat{\pi} = \pi_b$. In other words, only when $\gamma_i\left(\mathbf{p}^{\text{max}}\right) = \gamma_j\left(\mathbf{p}^{\text{max}}\right)$, for all $i, j \in \Omega$, the optimal uniform price can be analytically computed as $\pi_a = \pi^* = \hat{\pi} = \pi_b$. Otherwise, we only know $\pi_a \leq \pi^* \leq \hat{\pi} < \pi_b$, i.e., the optimal uniform price cannot be explicitly expressed in a closed form, although $\pi_a = \pi^* = \hat{\pi} < \pi_b$ may hold. If $\hat{\pi} = \pi_a$ holds, the optimal price $\pi^*$ is then clearly $\pi_a$. Based on this fact, we propose a low-complexity algorithm that gives the relay a sub-optimal price. Specifically, if we artificially increase $\pi_a$ and decrease $\hat{\pi}$ simultaneously until they meet at $\overline{\pi}$ and assume that

$$\rho(\pi) = \begin{cases} \pi \cdot \sum_{i=1}^{Q} \gamma_i(\mathbf{p}^{\text{max}}), & \text{if } 0 \leq \pi \leq \overline{\pi} \\ Q \cdot \left(\frac{1}{2} - \pi\right), & \text{if } \overline{\pi} < \pi \leq \frac{1}{2} \end{cases}, \qquad (2.17)$$

we can easily obtain the "optimal" uniform price as

$$\pi^* \approx \overline{\pi} = \frac{Q}{2\left[\sum_{i=1}^{Q} \gamma_i(\mathbf{p}^{\text{max}}) + Q\right]}. \qquad (2.18)$$

Generally speaking, setting (2.18) as the price can only result in a sub-optimal revenue for the relay. Nevertheless, the high computational complexity incurred by the exhaustive search is avoided and only limited information is needed to calculate (2.18): the number of active users in the network, i.e., $Q$, and the value of $\sum_{i=1}^{Q} \gamma_i(\mathbf{p}^{\text{max}})$. The relay can set a sufficiently low[10] price $\pi_s$, given which the NEP is $\mathbf{p}^{\text{max}}$, and find

---

[10]It can be verified that, if the relay sets a price $0 < \pi \leq \frac{1}{2} \min_{i \in \Omega} \left(1 + \gamma_i(p_i^{\text{max}}; \mathbf{0})\right)^{-1}$, then the resulting NEP of the user game is $\mathbf{p}^{\text{max}}$. The details are omitted due to the lack of space.

$\sum_{i=1}^{Q} \gamma_i(\mathbf{p}^{\max})$ by computing $\sum_{i=1}^{Q} \gamma_i(\mathbf{p}^{\max}) = \frac{\rho(\pi_s)}{\pi_s}$. The uniform price is determined in a similar way in the context of conventional cellular systems in [14] where the base station charges the users according to the throughput. Moreover, based on Corollary 1, we can establish the following corollary that guarantees the optimality of the uniform price in (2.18) when $Q = 1$.

**Corollary 2.** *When there is only one user in the network, the uniform price in (2.18) is the optimal one that generates higher revenues than any other uniform prices, i.e.,*

$$\overline{\pi} = \frac{1}{2\left[\gamma_1(p_1^{\max}) + 1\right]} = \arg \max_{\pi \geq 0} \left(\pi \cdot \gamma_1(p_1^*)\right). \tag{2.19}$$

*Furthermore, the transmit power of user $1$ is $p_1^{\max}$ at the NEP of the game $\mathcal{G}_{\mathrm{user}}$.* ■

We note that Corollary 2 directly follows Corollary 1 by invoking $\pi_a = \pi^* = \hat{\pi} = \pi_b$ when $Q = 1$. Furthermore, when there are sufficiently many users in the network or the users operate in low SINR regions, (2.18) is also a good approximation of the optimal uniform price. Specifically, when the number of users in the network is large, the sub-optimality of (2.18) can be explained as follows. It is natural that the level of interference observed by user $i$, i.e., $\sum_{j=1, j\neq i}^{Q} |g_j|^2 p_j$, increases when there are more active users. Hence, given a large value of $Q$, $\max_{i=1,2\cdots Q} \gamma_i(\mathbf{p}^{\max})$ becomes a small non-negative number due to the strong interference caused by the other users. Correspondingly, the difference between the lower bound and the upper bound on the optimal uniform price is not significant, i.e., $\hat{\pi} - \pi_a$ is a small number. Thus, the sub-optimal price (2.18), which lies between $\pi_a$ and $\hat{\pi}$, is close to the optimal one. Note that the small non-negative number $\hat{\pi} - \pi_a$ is also a upper bound on the gap between (2.18) and the optimal uniform price. Similar statements can be made when the network operates in low SINR regions as well. As in the existing literature (e.g., [14]), it is challenging to determine a priori the exact gap between (2.18) and the optimal uniform price, and hence, we shall verify in numerical results that the loss of

revenue is not significant in all the cases when the relay chooses (2.18), as compared to the optimal one obtained through exhaustive search, as its uniform price.

### 2.4.3 Differentiated Pricing With Complete Information

In this part, we extend the above analysis to a general case, in which different users may be charged at different prices, by considering that the relay has complete information about the network. It has been shown in [31] that the system performance can be improved if some users have complete information about the network. In the following analysis, the relay is assumed to know the maximum power constraints of all the users,[11] in addition to the channel coefficients. Under the differentiated pricing rule, we need to identify an optimal price vector $\mathbf{\Pi}^*$ set by the relay such that $\mathbf{\Pi}^* = \arg\max_{\mathbf{\Pi} \succeq \mathbf{0}} \left( \sum_{i=1}^{Q} \pi_i \gamma_i(p_i^*; p_{-i}^*) \right)$. Differentiated pricing is also referred to as price discrimination in the economics literature [33]. Similarly, depending on the channel conditions and maximum power constraints, the relay can charge different users at different prices. Before developing the differentiated pricing algorithm, we first express the optimal value of $\pi_i$ in terms of $\mathbf{p}^*$, for all $i \in \Omega$, in the following proposition.

**Proposition 1.** *Assume that $\mathbf{\Pi}^* = \{\pi_1^*, \pi_2^* \cdots \pi_Q^*\}$ is the optimal price vector, which generates the maximum revenue for the relay, and that $\tilde{\mathbf{p}}^* = \{\tilde{p}_1^*, \tilde{p}_2^* \cdots \tilde{p}_Q^*\}$ is the unique corresponding power allocation vector at the NEP of the user game $\mathcal{G}_{\mathrm{user}}$. Then, $\mathbf{\Pi}^*$ can be expressed in terms of $\tilde{\mathbf{p}}^*$ as follows*

$$\pi_i^* = \frac{1}{2(1 + \gamma_i(\tilde{\mathbf{p}}^*))}, \quad \forall i \in \Omega, \tag{2.20}$$

*where $\gamma_i(\tilde{\mathbf{p}}^*)$ is obtained by substituting $\tilde{\mathbf{p}}^*$ into (2.2).*

*Proof.* If $\mathbf{0} \prec \tilde{\mathbf{p}}^* \prec \mathbf{p}^{\mathrm{max}}$, then (2.20) directly follows the definition of NEP in (2.7)

---

[11]To implement the protocol, the user may be required to report its maximum transmit power level to the relay before entering the network.

and the first-order optimality condition. If $\tilde{p}_i^* = 0$ for some $i \in \Omega_0 \subset \Omega$, then $\pi_i^* \geq \frac{1}{2}$ and $\gamma_i(\tilde{\mathbf{p}}^*) = 0$ and hence, (2.20) also holds true. If $\tilde{p}_i^* = p_i^{\max}$ for some $i \in \Omega_{\max} \subset \Omega$, then $\tilde{\mathbf{p}}^*$ is the power allocation vector at the NEP of the game $\mathcal{G}_{\text{user}}$ for any $0 \leq \pi_i \leq \frac{1}{2(1+\gamma_i(\tilde{\mathbf{p}}^*))}$ and thus, it is clear that the optimal value of $\pi$ is $\pi_i^* = \frac{1}{2(1+\gamma_i(\tilde{\mathbf{p}}^*))}$. Therefore, Proposition 1 is proved. ∎

Proposition 1 enables us to express the price vector, which maximizes the relay's revenue, in terms of the transmit power levels at the NEP of the user-level game. Since the SINR is also a function of the transmit power levels of the users, we can then express the revenue, defined as the product of SINR and price, using a function of the power $\sum_{i=1}^{Q} \pi_i(\mathbf{p})\gamma_i(\mathbf{p})$. Therefore, instead of determining the optimal price vector directly, the relay can first decide the desired transmit power levels of the users and then set corresponding prices to enforce the users to transmit at these desired power levels. Mathematically, following Proposition 1 and substituting (2.20) into $\sum_{i=1}^{Q} \pi_i^* \gamma_i(p_i^*; p_{-i}^*)$, the problem of maximizing $\sum_{i=1}^{Q} \pi_i^* \gamma_i(p_i^*; p_{-i}^*)$ subject to $\mathbf{\Pi}^* \succeq \mathbf{0}$ can be reformulated as

$$\max \sum_{i=1}^{Q} \frac{\gamma_i(\tilde{\mathbf{p}}^*)}{2(1+\gamma_i(\tilde{\mathbf{p}}^*))} = \frac{\sum_{i=1}^{Q} \frac{|h_i|^2 p_{\mathcal{R}}}{|h_i|^2 p_{\mathcal{R}} + N_0} |g_i|^2 \tilde{p}_i^*}{2\left(\sum_{i=1}^{Q} |g_i|^2 \tilde{p}_i^* + N_0\right)} \tag{2.21}$$
$$s.t., \ \mathbf{0} \preceq \tilde{\mathbf{p}}^* \preceq \mathbf{p}^{\max},$$

where the objective function is linear-fractional and hence quasi-concave in $\tilde{\mathbf{p}}^*$ [113]. Therefore, the optimal value of $\tilde{\mathbf{p}}^*$ can be found by transforming (2.21) into a standard linear program [113], and the details of solving (2.21) are omitted due to the space limitations.

After the value of $\tilde{\mathbf{p}}^*$ is found, we can immediately obtain the optimal price vector using Proposition 1. It should be noted that, given the optimal price vector obtained using (2.20), the outcome of the game $\mathcal{G}_{\text{user}}$ when it reaches the NEP through iterations is that user $i$ transmits at the power of $\tilde{p}_i^*$, regardless of the initial power strategies. This

can be explained as follows. On one hand, we have shown in Proposition 1 that the optimal price vector $\boldsymbol{\Pi}^*$ can be expressed in (2.20) in terms of $\tilde{\mathbf{p}}^*$, i.e., one of the price vectors corresponding to $\tilde{\mathbf{p}}^*$ is $\boldsymbol{\Pi}^*$ given in (2.20). One the other hand, by uniqueness of the NEP of the game $\mathcal{G}_{\mathrm{user}}$ given any price vectors stated in Theorem 2, it can be seen that $\tilde{\mathbf{p}}^*$ is the unique NEP of the game $\mathcal{G}_{\mathrm{user}}$ if the relay sets $\boldsymbol{\Pi}^*$ as its pricing vector. Therefore, we can solve (2.21) to find $\tilde{\mathbf{p}}^*$ and then $\boldsymbol{\Pi}^*$ can be determined using (2.20). Furthermore, based on the objective function in (2.21), we have the following corollary regarding the upper bound on the revenue[12] of the relay.

**Corollary 3.** *The maximum revenue that the relay can obtain from all the users by charging the optimal differentiated prices is upper bounded by $\frac{1}{2}$, and for any $i \in \Omega$, $\lim_{|g_i|^2, |h_i|^2 \to \infty} \rho(\boldsymbol{\Pi}^*) = \frac{1}{2}$.* ∎

Corollary 3 states that, given differentiated prices, the maximum revenue of the relay can be collected from only one user if this user has a sufficient good channel condition. In other words, to maximize its revenue with complete information, the relay can set an appropriate price vector such that only one user transmits, if this user's channel gains are sufficiently large (i.e., $|g_i|^2, |h_i|^2 \to \infty$), while all the other users who are charged a price greater than or equal to $\frac{1}{2}$ remain silent. In contrast, under the uniform pricing algorithm, all the users are charged the same price according to (2.18) and hence, they will transmit simultaneously regardless of the channel conditions as long as the price is below $\frac{1}{2}$. Next, as a measure of comparison among different pricing schemes, we briefly discuss the average number of iterations played by the users and the information required by the relay to set the prices. Given complete information about the users, i.e., channel coefficients and power strategy space, the relay can directly compute the optimal differentiated price vector $\boldsymbol{\Pi}^*$, by solving the linear-fractional optimization problem in (2.21), and thus, it only needs to broadcast

---

[12]The unit of revenue is the same as that of the utility function, i.e., "nats/s/Hz" in this chapter [17]. Alternatively, the unit of the revenue can be converted to that of real money by multiplying the revenue with a constant converter without affecting the analysis [9].

Table 2.1: Average Number of Iterations and Information Requirement of Different Pricing Algorithms

| | Optimal Uniform | Sub-optimal Uniform | Differentiated |
|---|---|---|---|
| **Information Requirement** | $\varnothing$ | $Q = \vert\Omega\vert$ and $\sum_{i=1}^{Q} \gamma_i(\mathbf{p}^{\max})$ | $g_i, h_i$ and $\mathcal{P}_i = \left\{ p_i : 0 \le p_i \le p_i^{\max} \right\}$, for $i = 1, 2 \cdots Q$ |
| **Average Number of Iterations** | $m\bar{N}$ | $\bar{N}$ | 0 |

once the optimal price vector to the users. However, in the case of uniform pricing, the relay needs to set a sufficiently low price $\pi_s$ before identifying the sub-optimal uniform price, due to the constraint that only incomplete information about the users is available to the relay. Define $\bar{N}$ and $m$ as the average number of iterations required by the user game $\mathcal{G}_{\text{user}}$ to converge and the number of candidate quantized values of uniform prices, respectively. We list in Table 2.1 the average number of iterations prior to data transmissions of the users, and the information requirement of different pricing schemes.

### 2.4.4 System Utility Maximization

In the previous analysis, we have proposed two pricing mechanisms to maximize the relay's revenue, under the implicit assumption that the relay is solely revenue-driven. The proposed differentiated pricing algorithm, however, is also applicable if the relay wants to optimize the system utility which can be defined in any form. For the considered relay network, we have shown that, given any price vectors, there is a unique NEP in the user-level game, implying that the relay can set prices to enforce the users to transmit at desired power levels. Therefore, any system utility, defined as a function as the users' transmit power $\mathbf{p}$, can be achieved by setting appropriate prices.

As in [26], we denote the system utility which the relay wants to maximize as $U(\mathbf{p})$. Denote the optimal power levels maximizing $U(\mathbf{p})$ as $\bar{\mathbf{p}}$, i.e.,

$$\bar{\mathbf{p}} = \arg \max_{\mathbf{p} \in \mathcal{P}} U(\mathbf{p}). \tag{2.22}$$

33

After finding[13] $\bar{\mathbf{p}}$, the relay can set prices according to

$$\pi_i = \frac{1}{2(1 + \gamma_i(\bar{\mathbf{p}}))}, \quad \forall i \in \Omega. \tag{2.23}$$

Then, it is guaranteed that the users will transmit at $\bar{\mathbf{p}}$ at the NEP and thus, the system utility is maximized. For instance, let us take user scheduling as a concrete example. If the relay aims to schedule user 1 to transmit in a time slot and all the other users remain silent, it can set the price vector in such a way that $\pi_1 \in [0, \frac{1}{2})$ and $\pi_2 = \cdots = \pi_Q \geq \frac{1}{2}$ and, given this price vector, only user 1 will transmit when the game reaches the NEP. We state the pricing-based utility maximization problem formally in the following proposition.

**Proposition 2.** *Denote* $\bar{\mathbf{p}} = \arg\max_{\mathbf{p} \in \mathcal{P}} U(\mathbf{p})$, *where* $U(\mathbf{p})$ *is an arbitrary system utility. If the relay sets prices according to (2.23), then the system utility is maximized after the user-level game reaches the NEP.*

*Proof.* By plugging $\pi_i = \frac{1}{2(1 + \gamma_i(\bar{\mathbf{p}}))}$ and $p_{-i} = (\bar{p}_1 \cdots \bar{p}_{i-1}, \bar{p}_{i+1}, \cdots \bar{p}_Q)$ into the best response function of user $i$ given in (2.8), it can be shown that

$$\mathcal{B}_i(p_{-i}) = \bar{p}_i, \tag{2.24}$$

for $i = 1, 2 \cdots Q$. Therefore, $\bar{\mathbf{p}}$ is the NEP of the user game if the prices are set according to (2.23). Then, by uniqueness of NEP in $\mathcal{G}_{\text{user}}$, we see that $\bar{\mathbf{p}}$ must be the transmit power levels at the NEP corresponding to the price vector set based on (2.23). As a result, the system utility is maximized. ∎

Finally, we note that the relay is in fact taking the role of a *central planner* that has complete information about the network [26], if it wants to maximize the system

---

[13]If $U(\mathbf{p})$ is concave in $\mathbf{p}$, there exists efficient algorithms to maximize $U(\mathbf{p})$. Otherwise, the relay may need to maximize $U(\mathbf{p})$ via brutal-force search. As in [26], the details of optimizing $U(\mathbf{p})$ is beyond the scope of this chapter, wherein we focus on the design of pricing algorithms.

Figure 2.1: Sub-optimal Pricing and Convergence of Distributed Iterative Power Allocation Algorithms



Figure 2.2: Homogeneous Network: Average Revenue versus Average Channel Gain

Figure 2.3: Homogeneous Network: Average Sum Rate versus Average Channel Gain

utility which includes the revenue as a particular example. Nevertheless, the distinguishing feature of the proposed differentiated pricing algorithm is that it can enforce the users to transmit at desired power levels such that the system utility is maximized, even though these users are self-interested. Moreover, unlike in [26] wherein only near-optimal system utility can be achieved, we propose a pricing mechanism that can maximize any system utility by exploiting the uniqueness of NEP in the user-level game. The proposed pricing mechanism can be briefly described as follows. At the beginning of a frame, each user acquires its local information and, according to some performance metric (e.g., maximizing the revenue), the relay calculates the optimal power levels of all the users, sets its corresponding prices, and then announces the prices to the users. Then, Algorithm I is executed and the resulting NEP is achieved. In practice, Algorithm I can stop whenever the change in transmit power in two consecutive iterations is smaller than a certain threshold.

## 2.5 Numerical Results

For the convenience of illustration, $g_i$ and $h_i$ are modeled as independently Rayleigh distributed random variables, for $i \in \Omega$. The transmit power of the relay node and the maximum transmit power of each source node are normalized to one.

We consider a simple four-user network and randomly generate the channel gains and illustrate in Fig. 2.1 the convergence of the proposed distributed power allocation algorithm and the sub-optimal uniform pricing algorithm. The upper plot shows that the sub-optimal price (dashed line) is reasonably close to the optimal price (solid line) obtained through exhaustive search, which validates the use of (2.18) as the uniform price selected by the relay. Next, we evaluate and compare the proposed pricing algorithms based on two distinct performance metrics, i.e., average [14] revenue and sum rate.[15] An orthogonal transmission scheme (i.e., time-division multiple access, or TDMA, in this chapter) in which the source nodes do not interfere each other is also included in the comparison. Specifically, in the TDMA protocol, the users transmit in a round-robin manner and the relay charges each user using the optimal pricing scheme specified in (2.19).

### 2.5.1 Homogeneous Network Topology

Given the homogenous network topology, we assume that $g_i$ and $h_i$ have the same mean values, for $i \in \Omega$, i.e., $\mathbb{E}\{|g_1|^2\} = \mathbb{E}\{|h_1|^2\} = \cdots = \mathbb{E}\{|g_Q|^2\} = \mathbb{E}\{|h_Q|^2\}$, where $\mathbb{E}\{\cdot\}$ is the expectation operator.

---

[14]Throughout the simulations, "average" (e.g., average revenue, average rate) is taken over 10000 channel realizations.

[15]Due to the non-convexity, we solve the problem of sum rate maximization in (2.22) using greedy methods and obtain (locally) optimal solutions.
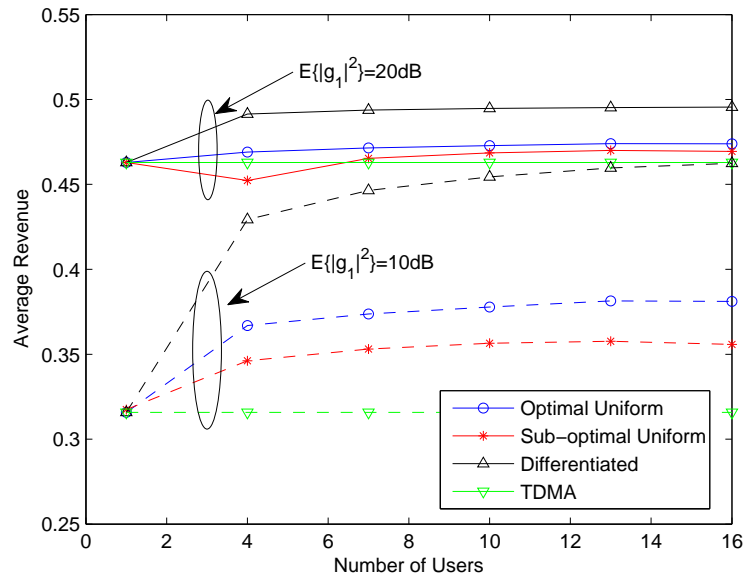
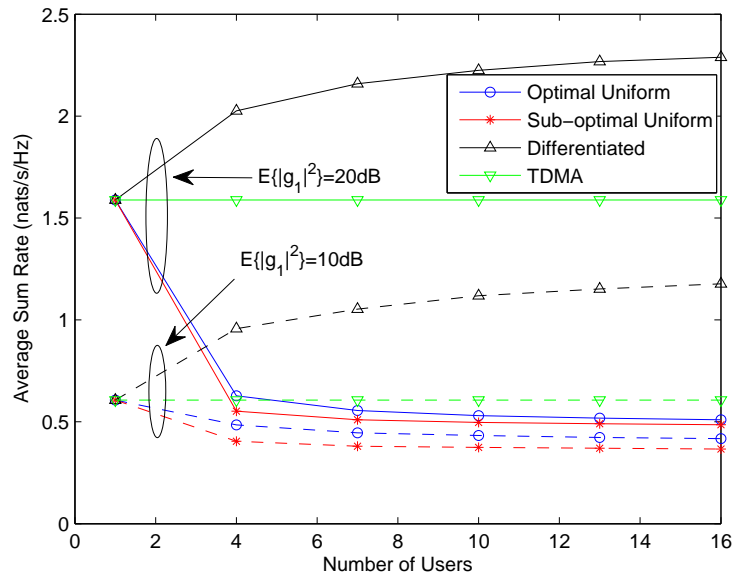Figure 2.4: Homogeneous Network: Average Revenue versus Number of Users



Figure 2.5: Homogeneous Network: Average Sum Rate versus Number of Users

### 2.5.1.1 Effects of Channel Gains

We consider a four-user network and examine the effects of channel gains on the average revenue and average sum rate in Fig. 2.2 and Fig. 2.3, respectively. As intuitively

expected and can be seen from (2.21), the average revenue of the relay increases as the channel condition becomes better. Fig. 2.2 demonstrates that the revenue loss due to the sub-optimality of the uniform price is not significant compared to the optimal uniform price. Among all the four pricing schemes, differentiated pricing generates the maximum revenue for the relay at the expense of having more information about the users. In other words, by allowing the users to transmit simultaneously, the differentiated pricing outperforms the optimal pricing in the TDMA protocol. This can be explained by noting that simultaneous transmission includes TDMA as a special case, i.e., simultaneous transmission reduces to TDMA if only one user is scheduled to transmit at a time (the other users can be charged a price greater than or equal to $\frac{1}{2}$ such that they transmit at a zero power). Regarding the upper bound on the revenues, it can be observed that the maximum revenue is always less than $\frac{1}{2}$ regardless of the channel conditions, which verifies Corollary 3. Fig. 2.3 compares the the proposed algorithms when they are applied to maximize the average sum rate of all the users (i.e., the system utility function becomes the sum rate). Note that, although the sub-optimal uniform pricing algorithm is applicable only for revenue maximization, we include the sub-optimal pricing for the completeness of comparison when we consider sum rate maximization. The optimal uniform price is numerically searched such that the sum rate is maximized. The proposed differentiated pricing achieves the highest average sum rate among all the considered protocols, since it can enforce the selfish users to transmit at the optimal power through pricing. For instance, if a user has a poor channel condition, the relay can charge this user a price greater than or equal to $\frac{1}{2}$ such that this user keeps silent and does not cause interference to the other users. We also observe from Fig. 2.3 that, when the channel condition is good enough, the TDMA protocol outperforms the two uniform pricing schemes, in which all the source nodes always transmit simultaneously and the heavy interference among the source nodes significantly limits the achievable rate.

Figure 2.6: Heterogeneous Network: Average Revenue versus Average Channel Gain



Figure 2.7: Heterogeneous Network: Average Sum Rate versus Average Channel Gain

### 2.5.1.2 Effects of Number of Users

In Fig. 2.4 and Fig. 2.5, we fix the average channel gains and vary the number of active users. Fig. 2.4 shows that, when there are more users competing for the relay,

40

the proposed pricing schemes achieve a higher revenue while the revenue obtained under the TDMA protocol does not change for the considered homogeneous network topology (since all the users with the same average channel statistics can be considered as one user in the TDMA protocol). Fig. 2.4 also indicates that the sub-optimal revenue of the relay gained by setting (2.18) as the uniform price is close to the optimal uniform one obtained through exhaustive search. Like in Fig. 2.2, the differentiated pricing outperforms its uniform counterpart and the TDMA protocol in terms of the average revenue. In terms of the average sum rate, the differentiated pricing is still able to achieve the best performance, and the revenue under the two uniform pricing schemes decreases when there are more users simultaneously transmitting in the network due to the strong interference.

### 2.5.2 Heterogeneous Network Topology

For the convenience of illustration, we assume that $\mathbb{E}\{|g_1|^2\} = \mathbb{E}\{|h_1|^2\}$ and $\mathbb{E}\{|g_2|^2\} = \mathbb{E}\{|h_2|^2\} = \cdots = \mathbb{E}\{|g_Q|^2\} = \mathbb{E}\{|h_Q|^2\}$ in heterogenous network topologies.

### 2.5.2.1 Effects of Channel Gains

As an example, we focus on a four-user network with heterogenous channel conditions in Fig. 2.6 and Fig. 2.7. Fig. 2.6 demonstrates that, among all the four pricing schemes considered in this chapter, the differentiated pricing yields the highest revenue for the relay, which is upper bounded by $\frac{1}{2}$. From Fig. 2.7, it can be seen that when the channel conditions becomes better, the average sum rate under the uniform pricing schemes are outperformed by that in the TDMA protocol and may not necessarily increase, since the interference also becomes stronger and reduces the received SINR. As in a homogeneous network topology, the proposed differentiated pricing achieves the highest average sum rate, since sum rate is only an instance of the system utility
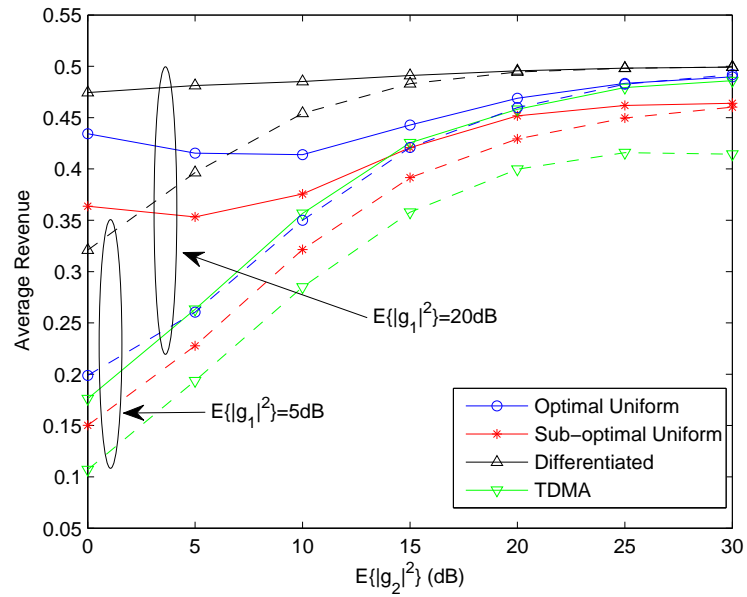
Figure 2.8: Heterogeneous Network: Average Revenue versus Number of Users



Figure 2.9: Heterogeneous Network: Average Sum Rate versus Number of Users

function and hence can be optimized using the differentiated pricing algorithm.

### 2.5.2.2 Effects of Number of Users

We consider fixed average channel gains and vary the number of users in a heterogeneous network topology in Fig. 2.8 and Fig. 2.9. It can be observed from Fig. 2.6 that the two uniform pricing schemes outperform the TDMA protocol in terms of the average revenue, and the proposed differentiated pricing achieves the highest revenue when there are more than one users regardless of the channel conditions. Fig. 2.9 shows that the differentiated pricing results in the highest average sum rate. Moreover, when the channel gains are strong, the average sum rate in the TDMA protocol is higher than that in the two uniform pricing schemes, since the strong interference can be avoided in the TDMA protocol.

To sum up, the proposed differentiate pricing scheme achieves the best performance in terms of the average revenue and sum rate, among all the four considered protocols. Compared to the TDMA protocol, the simultaneous transmission with uniform pricing schemes are generally more efficient in terms of the revenue, and less efficient in terms of the average sum rate (due to the unavoidable interference) when the channel gains are strong. Prior to concluding this part, we note that the analysis of general system utility maximization via the differentiated pricing in Chapter 2.4 is also valid and can be applied to arbitrary utility functions, though we do not show it in the simulations due to space limitations.

## 2.6 Conclusion

In this chapter, we considered a wireless relay network consisting of one relay node and multiple source-destination pairs. First, the interactions between the relay and the users were appropriately captured. We then modeled each user as a self-interested player, which aims at maximizing its own benefit by choosing the optimal transmit power, and analyzed the competition among the users using the notion of non-cooperative game

theory. It was proved that, in the non-cooperative game played by the users, there always exists a unique steady operating point, i.e., NEP, which can be achieved in a distributed manner. Next, under the assumption that the relay has only incomplete information about the users, we proposed a low-complexity algorithm, in which the relay charges the users at a sub-optimal uniform price. The analysis was then extended to differentiated pricing wherein the relay charges different users at different prices. We also showed that the proposed differentiated pricing can be applied to maximize any system utility. Extensive simulations were showed that the relay can gain the maximum revenue and the maximum sum rate by adopting the differentiated pricing algorithm, which though requires complete information about the users. Furthermore, given only incomplete information about the users, the relay can apply the proposed low-complexity sub-optimal uniform pricing algorithm without incurring a significant revenue loss as compared to the optimal uniform pricing algorithm. Interference cannot be avoided when using uniform pricing schemes and thus, the resulting average sum rate is less than that achieved by orthogonal transmission (e.g., TDMA) when the channels are in a good condition.

# CHAPTER 3

# Data Demand Dynamics in Wireless Communications Markets

In this chapter, we focus on the users' aggregate data demand dynamics in a wireless communications market served by a monopolistic wireless service provider (WSP). Based on the equilibrium data demand, we optimize the WSP's data plans and long-term network capacity decisions to maximize its profit. First, by considering a market where only one data plan is offered, we show that there exists a unique equilibrium in the data demand dynamics regardless of the data plans, and that the convergence of data demand dynamics is subject to the network congestion cost, which is closely related to the WSP's long-term capacity decision. A sufficient condition on the network congestion cost indicates that the WSP needs to provide a sufficiently large network capacity to guarantee the convergence of data demand dynamics. We also propose a heuristic algorithm that progressively optimizes the WSP's data plan to maximize its equilibrium revenue. Next, we turn to a market where two different data plans are offered. It is shown that the existence of a unique equilibrium data demand depends on the data plans, and the convergence of data demand dynamics is still subject to the network congestion cost (and hence, the WSP's network capacity, too). We formalize the problem of optimizing the WSP's data plans and network capacities to maximize its profit. Finally, we discuss the scenario in which the data plans are offered by two competing WSPs and conduct extensive simulations to validate our analysis.

## 3.1 Introduction

We have witnessed over the last decade a successful proliferation of wireless networks, which support a variety of services and applications, and increasingly heated competition among the wireless service providers (WSPs). To sustain their competitive positions in the market and increase revenues, WSPs themselves will need to appropriately price their scarce network resources and expand their network capacities to support the unprecedented amount of wireless traffic. Hence, it becomes of paramount importance for these WSPs to understand how the aggregate data demand of all the subscribers evolves and how the demand is affected by various pricing plans.

In this chapter, we are interested studying the users' aggregate data demand dynamics, and optimizing the WSP's data plans and network capacities in a wireless communications market. In general, the WSP's network capacity is difficult to change once it is deployed and hence, it is a long-term strategy for the WSP to decide its network capacity [41]. In contrast, the WSP can adjust its data plans over the lifespan of its network infrastructure, although the data plans cannot be updated as frequently as the users change their data plan subscription. Overall, we will assume that the users may change their plan frequently, based on their short-term (e.g., a few days or weeks per period) decisions, the WSP's data plans are changed less frequently, based on the WSP's medium-term (e.g., several months or years per period) decisions, while the WSP's network capacity decision is a long-term (e.g., several years per period) decision. In order to evaluate and compare the long-term profitability of networks with different capacities, the WSP needs to predict its maximum profit for each network capacity configuration. To maximize revenue given the network capacity and the associated cost, the WSP needs to know the users' aggregate data demand and their willingness to pay for the service, and then choose its optimal data plans. Hence, by using backward induction, we study first the users' dynamic decisions as to whether or not they subscribe to the WSP's data plans (i.e., short-term problem), then the WSP's

revenue-maximizing data plans (i.e., medium-term problem), and finally the WSP's network capacity decision (i.e., long-term problem). Note that we assume in our study that the medium-term period is sufficiently short compared to the long-term period, while it is sufficiently long compared to the short-term period.

We consider a wireless market with a monopolistic WSP serving a sufficiently large number of users. For the sake of analysis, we consider that the WSP can offer one or two data plans, while each user can subscribe to one of the available data plans. Due to the resource constraint (e.g., network capacity), congestion effects are observed when multiple users share the same network, degrading the network performance (e.g., increasing delays). Essentially, congestion effects are a type of negative network externalities and have similar impacts to prices on the users' experiences (i.e., utilities). Thus, congestion effects are also referred to as *congestion costs* in the literature [46][55]. Taking into consideration the charged price and congestion cost, each user can dynamically decide whether to subscribe to the WSP's service and which data plan to subscribe to. First, by considering a market where only one data plan is offered, we show that there exists a unique equilibrium in the data demand dynamics regardless of the data plan or congestion costs. Nevertheless, the convergence of data demand dynamics is subject to the network congestion cost, which is closely related to the WSP's long-term capacity decision. We derive a sufficient condition for the convergence of data demand dynamics, indicating that that the WSP needs to provide a sufficiently large network capacity. A heuristic algorithm is also proposed to progressively optimize the WSP's data plan such that its equilibrium revenue is maximized. Next, we turn to a market where two different data plans are offered. We show that the existence of a unique equilibrium data demand depends on the data plans. Moreover, the convergence of data demand dynamics is still subject to the network congestion cost (and hence, the WSP's network capacity, too). The problem of optimizing the WSP's data plans and network capacities is formalized and solved by numerical meth-

ods to maximize its profit. Next, we discuss the scenario in which the data plans are offered by two competing WSPs (i.e., a duopoly market) and find that the two WSPs only need to adjust their data plans a few times before reaching an equilibrium. Finally, extensive simulations are conducted to validate our analysis. Numerical results shows that, to maximize the profit, the WSP needs to increase the network capacity for its capped data plan while reducing the network capacity for its unlimited data plan. This coincides with the current trend that some WSPs have discontinued the offering of unlimited data plans [52].

The rest of this chapter is organized as follows. We review the related literature in Chapter 3.2. Chapter 3.3 describes the model. In Chapter 3.4 and Chapter 3.5, we study the data demand dynamics, data plan decision and capacity decision for a wireless market where one and two data plans are offered, respectively. In Chapter 3.6, we provide numerical results to validate our analysis. Finally, we conclude this chapter in Chapter 3.7.

## 3.2 Related Works

The engineering community has recently started to analyze as well as consider the design of existing and emerging wireless markets from various perspectives. Because of the space limitation, we only provide an incomplete list of related literature. In our previous work [35], we study the user subscription dynamics and revenue maximization in both monopoly and duopoly communications markets, based on a general distribution of users' valuation of quality-of-service (QoS) and a general QoS function that captures negative network externalities. Focusing on two specific access technologies (i.e., wide and local area network), the authors in [36] apply a stochastic geometric model and study the convergence of user subscription dynamics. In [37], the authors showed that non-cooperative communications markets suffer from unfair revenue dis-

tribution among the service providers and proposed a revenue-sharing mechanism that requires cooperation among the service providers. The behavior of users and its impact on the revenue distribution, however, were not explicitly considered in [37]. [66] studies technology adoption and competition between incumbent and emerging network technologies. The model characterizing the users' valuation of QoS is restricted to uniform distributions, and only constant QoS functions and positive network effects are considered in [66]. The user evolution in wireless social community networks is investigated in [38], where a key assumption is that the social community network provides a higher QoS to each user as the number of subscribers increases. While this assumption is valid if the network coverage is the only factor that determines the QoS, it does not model the QoS degradation due to, for instance, user traffic congestions at the WSP. By taking into account the congestion cost (i.e., negative network externality), [39] studies the feasibility of Paris Metro pricing (PMP) and shows sufficient conditions on the congestion cost functions, under which PMP leads to a higher revenue or social welfare than flat-rate pricing. Pricing decisions (restricted to unlimited data plans) and network capacity decisions in the presence of network congestion effects are studied in [41], where a missing part is the analysis of users' subscription decisions. [43] investigated market dynamics emerging when next-generation networks and conventional networks coexist, by applying a market model that consists of content providers, service providers, and users. Nevertheless, the level of QoS that a certain technology can provide was not considered in the model. The authors in [99] formulate a rate allocation problem by incorporating the participation of content providers into the model, and derive equilibrium prices and data rates. In [45], time-dependent pricing is studied from the perspective of its efficiency in terms of revenues. In [46], an upper bound on the efficiency loss as a result of price competition is derived in the context of congested markets, where an infinite number of users can selfishly route their traffic through the network.

In the aforementioned works, however, several key points are neglected. First, user heterogeneity in terms of data demand is not considered in these works (except for [99]). Specifically, it has been an implicit yet common assumption in these works that every user has the same data demand when it subscribes to the service provider. In other words, the QoS provided by the service provider only depends on the number of subscribers, regardless of their actual demand. Hence, user heterogeneity in terms of data demand cannot be captured and the QoS characterization may not be accurate under this assumption. Second, in most of the works, only a single data plan (e.g., flat-rate or "unlimited", usage-based price) is considered. Nevertheless, with the exploding popularity of smart phones, multiple pricing schemes are emerging in the market. For instance, capped data pricing plans and unlimited data pricing plans are both available in current wireless markets. Last but not least, it remains unknown how the congestion costs affect the aggregate data demand dynamics, in terms of both the equilibrium point and convergence, and the resulting revenue of the WSP. To address all these concerns, we propose a unified model that captures the user heterogeneity in terms of data demand and various practical data plans. Then, we study the users' data demand dynamics, and the WSP's data plan decision and network capacity decision.

## 3.3  Model

Consider a wireless communications market where one monopolistic WSP, denoted by $\mathcal{W}$, offers to $N$ users data communications service, which takes up an overwhelming majority of the wireless traffic [53]. By assuming that $N$ is sufficiently large such that each user is negligible,[1] we use a continuum user population model and normalize the number of users to $1$ [35]–[41]. In general, WSP $\mathcal{W}$ may offer multiple data plans, and users can choose any of the plans depending on their own preferences (the user

---

[1]Another interpretation of the continuum model is that there is a representative user which has the same characteristics (e.g., data demand) as each user $i$ in the market with a certain probability.

choice shall be detailed later). As in [41], to keep the analysis tractable, we assume that WSP $\mathcal{W}$ offers up to two data plans, represented by $\mathcal{P}_1$ and $\mathcal{P}_2$, respectively. For notational convenience, we also refer to users that subscribe to the plan $\mathcal{P}_i$ as $\mathcal{P}_i-$users (or $\mathcal{P}_i-$subscribers), for $i = 1, 2$. Next, we shall provide the modeling details of the WSP and users.

### 3.3.1  WSP Model

Before entering a market, the WSP needs to first make investment in infrastructure. In this chapter, we concentrate on the WSP' capacity deployment which, once determined, is difficult to adjust and hence is an irreversible long-term decision [40]. Denote by $C_i \geq 0$ the network capacity (normalized by the number of users $N$) that the WSP allocates to its data pricing plan $\mathcal{P}_i$, for $i = 1, 2$. Assuming that the WSP incurs an average cost of $\tau$ per unit capacity,[2] we can express the WSP's equilibrium profit per short-term period (i.e., users' subscription period) as

$$\Pi_{\mathcal{W}} = \sum_{i=1,2} \{R_i - \tau C_i\}, \qquad (3.1)$$

where $R_i$ is the equilibrium revenue per short-term period derived from $\mathcal{P}_i-$users. Note that in (3.1), we neglect the recurring cost of serving the users, which can also be absorbed into the revenue $R_i$ [40]. To maximize its profit given the users' rational decisions, the WSP shall strategically determine its capacity $\mathbf{C} = \{C_1, C_2\}$. After building the network, the WSP decides its data plans and may alter them throughout the network's lifespan.

In today's wireless market, the most popular data plans are "unlimited", "capped" and "usage-based", all of which can be represented by a unified pricing model specified

---

[2]The cost is averaged over the lifespan of the network infrastructure. For instance, if a network with a lifespan of $T$ short-term periods (i.e., users' subscription period) is built at a cost of $\tilde{\tau}$ per unit capacity, then the average cost per unit capacity is $\tau = \tilde{\tau}/T$.

by $(p, d^*, \gamma)$: each subscriber pays a fixed subscription fee $p$ that allows it to transmit and receive up to $d^*$ units of data; for each unit of additional data usage exceeding the capped data limit $d^*$, the subscriber pays $\gamma$. In special cases, a capped data plan characterized by $(p, d^*, \gamma)$ becomes a usage-based one if $p = 0$ and $d^* = 0$, and an unlimited data plan if $d^* = \infty$ or $\gamma = 0$. For analytical tractability and to gain insights on how the congestion costs affect the data demand dynamics, we assume that the WSP's data plan $\mathcal{P}_1 = (p_1, +\infty, 0)$ is "unlimited" whereas its data plan $\mathcal{P}_2 = (p_2, d_2^*, \gamma_2)$ is "capped".[3] This assumption, which may seem strong, can be justified by noting that some WSPs have (partially) resorted to capped data plans in view of the soaring wireless data service demand that frequently clogs their network infrastructure.[4] Moreover, even if the WSP offers two capped data plans, it is likely that one of the data plans has a very high data limit, which only a negligible fraction of subscribers can exceed in practice, and thus this data plan is almost "unlimited" (see, e.g., [52]).

### 3.3.2  User Model

Due to the capacity constraint, the network becomes more congested (i.e., negative network externalities or effect) as more data flow is transmitted [36][45]. Such an effect is quantified by the congestion cost, which has similar impacts to prices on the users' experiences (i.e., utilities) [46]. We denote the congestion cost associated with the data pricing plan $\mathcal{P}_i$ by $g_i(D_i, C_i)$, where $i = 1, 2$ and $D_i \geq 0$ is the aggregate data demand (i.e., the total data demand of all the $\mathcal{P}_i$-users over a certain period) and $C_i$ is the capacity allocated to $\mathcal{P}_i-$users. Without causing ambiguity, we simplify $g_i(D_i, C_i)$ as $g_i(D_i)$ by removing $C_i$ wherever applicable. An implicit assumption in the model is that congestion costs for different data plans are independent of each other, which

---

[3]In the most general case where both data plans are "capped", the approach of analysis in this chapter is still applicable, although the analysis becomes more complicated.

[4]Starting from June 7, 2010, AT&T discontinued unlimited data plans to its new iPhone users and adopts a capped data plan as considered in this chapter [52].

may be achieved by splitting network capacity among the plans [40].

Users are heterogeneous in the sense that they may have different data service demand and different benefits of utilizing the WSP's communications service. To model the user heterogeneity, each user $k$ is characterized by a two-element tuple $(\theta_k, d_k)$, where $\theta_k$ indicates user $k$'s benefit from data service and $d_k$ denotes its data demand over a certain period (e.g., a month or a day). The values of $\theta_k$ and $d_k$ can be determined by various approaches. For instance, $(\theta_k, d_k)$ may be user $k$'s intrinsic characteristic and not influenced by the WSP's pricing schemes. In such scenarios, each individual user has *inelastic* demand [45][46], although the aggregate demand of all the users is still elastic and influenced by the prices. Mathematically speaking, when user $k$ subscribes to the WSP's data plan $\mathcal{P}_i$, its utility is given by

$$u_{k,i} = \theta_k - g_i(D_i) - p_i - \gamma_i[d_k - d_i^*]^+, \tag{3.2}$$

where $[\,x\,]^+ = \max\{0,\,x\}$, and if its data demand exceeds the granted data limit $d_i^*$, the term $\gamma_i[d_k - d_i^*]^+$ is positive and represents the additional cost user $k$ incurs. Similar utility functions have been used in [66][39][99][44][45] and references therein. The utility function in (3.2) can be interpreted as follows: $\theta_k$ represents the benefit that user $k$ receives from $d_k$ units of data service, $g_i(D_i)$ indicates the congestion cost (i.e., negative network externality), and $p_i + \gamma_i[d_k - d_i^*]^+$ is the payment made to WSP $W_i$ [44]. Users that do not subscribe to any data plans obtain zero utility. Now, we impose some standard assumptions on the users' data demand and their benefits, users' subscription decisions, and the congestion function $g_i(D_i)$.

*Assumption 1:* The users' benefits and their data demand follow a two-dimensional distribution whose joint density function $f(\theta, d)$ is defined on $\mathcal{U} = \{(\theta, d)\,|\,0 \leq \theta \leq \theta_{\max}, 0 \leq d \leq d_{\max}\}$. For completeness of definition, we have $f(\theta, d) = 0$ for all $(\theta, d) \notin \mathcal{U}$. The cumulative density function is given by $F(\theta, d) = \int_{-\infty}^{d} \int_{-\infty}^{\theta} f(x, y)dxdy$

for $(\theta, d) \in \mathbb{R}^2$.

*Assumption 2:* Each user $k$ subscribes to the data plan $\mathcal{P}_i$ if $u_{k,i} > u_{k,j}$ and $u_{k,i} \geq 0$ for $i, j \in \{1, 2\}$ and $i \neq j$. If $u_{k,1} = u_{k,2} \geq 0$, user $k$ subscribes to the unlimited data plan $\mathcal{P}_1$.[5]

*Assumption 3:* $g_i(D_i)$ is a non-negative, non-decreasing and differentiable[6] function in $D_i \in [0, D_{\max}]$, where $D_{\max}$ is the maximum possible aggregated data demand, normalized with respect to the total population, and given by

$$D_{\max} \triangleq \int_{y=0}^{d_{\max}} \int_{x=0}^{\theta_{\max}} y f(x, y) dx dy. \tag{3.3}$$

We briefly explain the above three assumptions. Assumption 1 can be considered as an expression of user diversity in terms of the benefits and their data demand. The lower bound on the interval is set as zero to simplify the analysis, and this will be the case when there is enough diversity in the users so that there are non-subscribers for any positive price [40][41]. Assumption 2 captures the user rationality. A rational user will subscribe to the data plan that provides a higher utility if at least one data plan provides a non-negative utility, and to neither data plan otherwise. Assumption 3 indicates an intuitive fact that the congestion cost that each user experiences when subscribing to the data plan $\mathcal{P}_i$ becomes larger when the aggregate data demand increases.

Before concluding this part, it is worthwhile to provide the following remarks regarding our model.

*Remark 1:* As in [45], we assume for the convenience of analysis that each individual user $k$ has an inelastic and fixed demand $d_k$ (and benefit $\theta_k$, too). Alternatively, $d_k$ can be determined by solving a utility maximization problem and $\theta_k$ is the maximum

---

[5]Online surveys show that users generally prefer an unlimited data plan to a capped one [54]. Moreover, specifying an alternative tie-breaking rule (e.g., random selection between the two data plans) in case of $u_{k,1} = u_{k,2} \geq 0$ will not significantly affect the analysis of this chapter.

[6]Since $g_i(\cdot)$ is defined on $[0, D_{\max}]$, we use a one-sided limit to define the derivative of $g(\cdot)$ at 0 and $D_{\max}$, e.g., $g_i'(0) = \lim_{D_i \to 0^+} [g_i(D) - g_i(0)]/(D_i - 0)$.

benefit that user $k$ receives [99]. Nevertheless, given the WSP's data plans, $(\theta, d)$ still follows a certain distribution over all the users and thus, our approach can be viewed as a proxy to determine the users' demand and benefit, provided that the distribution does not change significantly with the data plans.

*Remark 2:* Compared to the congestion cost function used in the existing literature that disregards the user heterogeneity in terms of data demand and is defined solely in terms of the number of subscribers [35][39][41], $g_i(D_i)$ is more accurate in modeling the congestion effect. Whilst the actual congestion cost also depends on when the users utilize the network, we consider the congestion cost *averaged* over time and ignore the time dependency to keep the analysis tractable [46].

*Remark 3:* The shape of the congestion cost function $g_i(D_i)$ may be determined by various factors, including the network capacity, resource allocation schemes and/or scheduling algorithms used for the data plan $\mathcal{P}_i$. While our analysis applies to a general function $g_i(D_i)$ satisfying Assumption 3, we shall explicitly focus on the impacts of network capacities on $g_i(D_i)$ when we derive specific results or study the WSP's long-term capacity decision. For instance, a concrete example is given by $g_i(D_i) = D_i/C_i$, which has been widely used (with minor modification, e.g., assuming all the users have the same data demand) in the prior work [39][41][50][51].[7]

*Remark 4:* In addition to negative network externalities (i.e., congestion costs in this chapter), positive network externalities may also be observed in a communications network. For instance, when more users subscribe to the WSP's data plan, the value of communications service may become higher as more users can communicate with each other [66]. As in prior research (e.g., [36][39][40][41][46][47][55]), we neglect the positive network externalities and concentrate on the impacts of congestion effects on the users' subscription decisions.

---

[7]Another congestion cost function widely adopted in the literature is $g_i(D_i) = 1/(C_i - D_i)$, which satisfies Assumption 3. Thus, our analysis is also applicable if $g_i(D_i) = 1/(C_i - D_i)$ is considered.

## 3.4 Wireless Communications Market: Single Data Plan

In this part, we study the wireless communications market where the WSP offers a single data plan. Without loss of generality, we assume that the offered data plan is $\mathcal{P}_2 = (p_2, d_2^*, \gamma_2)$ and, as aforementioned, the unlimited data plan $\mathcal{P}_1 = (p_1, \infty, 0)$ is a special case when the data limit is infinity. The timing (i.e., order of moves) can be described as follows.

**Stage 1 (long-term):** The WSP decides its network capacity $C_2$ to deploy to maximize its profit.

**Stage 2 (medium-term):** Given $C_2$, the WSP chooses its optimal data plan $\mathcal{P}_2 = (p_2, d_2^*, \gamma_2)$ by specifying $p_2$, $d_2^*$ and $\gamma_2$ to maximize its revenue.

**Stage 3 (short-term):** By jointly considering the congestion cost and offered data plan, users decide whether or not to subscribe to the WSP's service.

From the described timing, we see that the WSP can be regarded as the leader whereas the users are followers. Thus, in order to identify the optimal data plan and network capacity, the WSP needs to first know how the users make their subscription decisions. Therefore, we proceed with our analysis using backward induction.

### 3.4.1 Users' Subscription Decisions

Due to rationality, users will not choose to subscribe to the WSP's data plan $\mathcal{P}_2$ if they cannot obtain non-negative utilities. Essentially, the subscription decision stage can be formalized as a non-cooperative game with an infinite number of players, the solution to which is (Nash) equilibrium. At an equilibrium, if any, no users can gain more benefits by deviating from their decisions. In other words, the aggregate data demand of those users subscribing to the WSP's data plan does not change at the equilibrium. Given the WSP's long-term capacity decision and an aggregate data demand $D_2$ of the subscribers, the congestion cost is uniquely given by $g_2(D_2)$. Moreover, the

users' subscription decisions are also determined based on the sign of (3.2), i.e., user $k$ subscribes to the data plan $\mathcal{P}_2$ if and only if $u_{k,2} = \theta_k - g_2(D_2^*) - \gamma_2[d_k - d_2^*]^+ \geq 0$. Hence, we study the users' subscription decisions at the equilibrium by specifying the equilibrium (aggregate) data demand $D_2^*$. First, we can mathematically express the equilibrium data demand as

$$D_2^* = h_2(D_2^*) = \int_{y=0}^{d_{\max}} \int_{x=g_2(D_2^*)+p_2+\gamma_2[y-d_2^*]^+}^{\theta_{\max}} yf(x,y)dxdy. \tag{3.4}$$

Thus, an equilibrium data demand exists if and only if the mapping $h_2(D_2^*)$ in (3.4) has at least one fixed point. Next, we formally define the equilibrium data demand as follows.

*Definition 1:* When only data plan $\mathcal{P}_2$ is offered, $D_2^* \in [0, D_{\max}]$ is an *equilibrium* data demand if it satisfies

$$h_2(D_2^*) = D_2^*. \tag{3.5}$$

We establish in the following proposition the existence and uniqueness of an equilibrium data demand $D_2^*$.

**Proposition 3.** *For any data plan $\mathcal{P}_2 = (p_2, d_2^*, \gamma_2)$, there exists a unique equilibrium data demand satisfying* (3.5).

*Proof.* To facilitate the proof, we first define an auxiliary function $\tilde{h}_2(D_2) = h_2(D_2) - D_2$ for $D \in [0, \infty)$, where $h_2(\cdot)$ is defined in (3.7). By Definition 1, $D_2^*$ is an equilibrium point if and only if it is a root of $\tilde{h}_2(\cdot)$. Hence, it suffices to show that $\tilde{h}_2(\cdot)$ has a unique root on its domain.

Let $D_{\max} \geq D_{2,a} > D_{2,b} \geq 0$ be two arbitrarily-chosen real numbers. Then, it

follows that

$$\tilde{h}_2(D_{2,a}) - \tilde{h}_2(D_{2,b}) = h_2(D_{2,a}) - D_{2,a} - [h_2(D_{2,b}) - D_{2,b}]$$

$$= -\int_{y=0}^{d_{\max}} \int_{x=g_2(D_{2,b})+p_2+\gamma_2[y-d_2^*]^+}^{g_2(D_{2,a})+p_2+\gamma_2[y-d_2^*]^+} yf(x,y)dxdy - (D_{2,a} - D_{2,b}).$$

(3.6)

Since $g_2(\cdot)$ is non-decreasing in $[0, D_{\max}]$, we have $g_2(D_{2,a})+p_2+\gamma_2[y-d_2^*]^+$ is greater than or equal to $g_2(D_{2,b})+p_2+\gamma_2[y-d_2^*]^+$ and hence, $\int_{y=0}^{d_{\max}} \int_{x=g_2(D_{2,b})+p_2+\gamma_2[y-d_2^*]^+}^{g_2(D_{2,a})+p_2+\gamma_2[y-d_2^*]^+} yf(x,y)dxdy$ is non-negative. Thus, it can be seen that $\tilde{h}_2(D_{2,a}) - \tilde{h}_2(D_{2,b}) \leq -(D_{2,a} - D_{2,b}) < 0$ for any $D_{\max} \geq D_{2,a} > D_{2,b} \geq 0$. That is, the auxiliary function $\tilde{h}_2(\cdot)$ is strictly decreasing in $[0, D_{\max}]$.

On the one hand, $\tilde{h}_2(0) = h_2(0) - 0 \geq 0$, and on the other hand, $\tilde{h}_2(D_{\max}) = h_2(D_{\max}) - D_{\max} \leq 0$ Since $\tilde{h}_2(\cdot)$ is continuous on $[0, D_{\max}]$, we see that $\tilde{h}_2(\cdot)$ has a unique root $D_2^* \in [0, D_{\max}]$, by applying the intermediate value theorem. This proves Proposition 3. ∎ □

It can be seen from Proposition 3 that the data plan $\mathcal{P}_2 = (p_2, d_2^*, \gamma_2)$ uniquely determines the equilibrium data demand. Although it is in general rather difficult to express $D_2^*$ as an explicit function of $\mathcal{P}_2 = (p_2, d_2^*, \gamma_2)$, we summarize in Proposition 4 the relation between the data plan $\mathcal{P}_2 = (p_2, d_2^*, \gamma_2)$ and the equilibrium data demand $D_2^*$.

**Proposition 4.** *For any congestion cost function $g_2(D_2)$ satisfying Assumption 3, the equilibrium data demand $D_2^*$ has the following properties:*

*1. $D_2^* > 0$ if and only if $0 \leq p_2 < [\theta_{\max} - g_2(0)]^+$;*

*2. $D_2^*$ is non-increasing in $p_2 \geq 0$;*

*3. $D_2^*$ is non-increasing in $\gamma_2 \geq 0$;*

*4. $D_2^*$ is non-decreasing in $d_2^* \geq 0$.* ∎

Property 1 shows that no users will subscribe to the WSP's data plan $\mathcal{P}_2$ if the

fixed subscription fee $p_2$ exceeds the maximum benefit among all the users minus the minimum congestion cost. Properties 2 and 3 are consistent with the standard demand-price relation: increasing the price will not increase the demand. Property 4 indicates that the data demand will increase or at least remain the same if the data limit $d_2^*$, which each subscriber can enjoy without incurring additional costs, increases. This stems from the fact that increasing the data limit results in the decrease of payment for users with high data demand exceeding the limit $d_2^*$.

In practice, the users do not have complete information regarding each other and hence, they may not make directly the subscription decisions that lead to an equilibrium. Instead, an adjustment process where the users update their subscription decisions based on limited information is required. To formally describe the adjustment process, we consider a discrete-time model denoted by $\{D_2^t\}_{t=0}^{\infty}$, where $D_2^t \geq 0$ is the (aggregate) data demand in the $t$-th time period and $D_2^0 \in [0, D_{\max}]$ is the initial data demand. A natural and well-studied approach to modeling the adjustment process is the best-response dynamics, in which each decision maker chooses the best action in response to the decisions made by the others. As in [35][38][42][47], we consider the best-response dynamics based on naive (or static) expectation, and assume that the users can only change their subscription decisions (e.g., opt out of the plan $\mathcal{P}_2$) at discrete time periods indexed by $t = 1, 2, \cdots$. Specifically, at the beginning of the time period $t$, user $k$ holds a (static) belief on the congestion cost, denoted by $\tilde{g}_{2,k}(D_2^t) = g_2(D_2^{t-1})$, and makes its subscription decision in a myopic way [66][38].[8] When only one data plan $\mathcal{P}_2$ is offered by the WSP, each user has a choice of whether to subscribe to the plan at the beginning of each time period. In particular, user $k$ subscribes to the data plan $\mathcal{P}_2$ in the time period $t$ if and only if it believes that its utility $\theta_k - \tilde{g}_{2,k}(D_2^t) - p_2 - \gamma_2[d_k - d_2^*]^+ = \theta_k - g_2(D_2^{t-1}) - p_2 - \gamma_2[d_k - d_2^*]^+ \geq 0$. Note

---

[8]This model of belief formation is called naive or static expectations [56]. A similar dynamic model of belief formation and decision making has been extensively adopted in the existing literature such as [35][66][38][42].

that, in order to make subscription decisions at time $t$, the users need to know the data plan $\mathcal{P}_2$ and receive a signal indicating the congestion cost $g_2(D_2^{t-1})$ at $t-1$. The best-response decision model implies that, for $t = 1, 2, \ldots$, the data demand dynamics evolves following a sequence $\{D_2^t\}_{t=0}^{\infty}$ specified by

$$D_2^t = h_2(D_2^{t-1}) = \int_{y=0}^{d_{\max}} \int_{x=g_2(D_2^{t-1})+p_2+\gamma_2[y-d_2^*]^+}^{\theta_{\max}} yf(x,y)dxdy, \qquad (3.7)$$

starting from an initial point $D_2^0 \in [0, D_{\max}]$. Essentially, the dynamics in (3.7) is a fixed point iteration for $h_2(\cdot)$ and it converges regardless of the initial point if $|h_2'(D_2)| < 1$ for $D_2 \in [0, D_{\max}]$ [69]. Nevertheless, $|h_2'(D_2)| < 1$ may not hold for all congestion cost functions, resulting in oscillation in the data demand dynamics. In accordance, the WSP's revenue becomes instable and may causes higher risks for the WSP's operation in the market. Let us consider a hypothetical example to explain this point. Suppose that the network is highly underutilized in the time period $t$ and each subscriber incurs a low congestion cost. Users expect that the congestion cost will remain low in the period $t + 1$, and thus more users subscribe to the data plan $\mathcal{P}_2$, leading to a high congestion cost in the time period $t + 1$. The increase of congestion cost in turn will induce a small amount of data demand in the time period $t + 2$. When the congestion cost function is very sensitive to the aggregate data demand, the data demand dynamics may oscillate around or diverge away from the equilibrium point. In the following proposition, we provide a sufficient condition under which the data demand dynamics is guaranteed to converge regardless of the initial points.

**Proposition 5.** *For any data plan $\mathcal{P}_2 = (p_2, d_2^*, \gamma_2)$, the data demand dynamics specified by (3.7) converges to the unique equilibrium point starting from any initial point $D_2^0 \in [0, \bar{d}]$ if*

$$\max_{D_2 \in [0, D_{\max}]} g_2'(D_2) < \frac{2}{K \cdot d_{\max}^2}, \qquad (3.8)$$

*where $D_{\max}$ is given by (3.3), $d_{max}$ is the maximum individual demand and $K =$*

$\max_{(\theta,d)\in\mathcal{U}} f(\theta,d)$.

*Proof.* We prove the convergence of the data demand dynamics based on contraction mapping theorem, which is formally stated as follows [57].

*Definition 2 [57]:* A mapping $\mathbf{T} : \mathcal{X} \to \mathcal{X}$, where $\mathcal{X}$ is a closed subset of $\mathbb{R}^n$, is called a contraction if there is a real number $\kappa \in [0,1)$ such that

$$\left\| \mathbf{T}(x_1) - \mathbf{T}(x_2) \right\| \leq \kappa \cdot \left\| x_1 - x_2 \right\|, \quad \forall\, x_1, x_2 \in \mathcal{X}, \tag{3.9}$$

where $\|\cdot\|$ is some norm defined on $\mathcal{X}$.

Proposition 1.1 in Chapter 3 of [57] shows an important property of a contraction mapping $\mathbf{T}$ that the update sequence generated by $x^{t+1} = \mathbf{T}(x^t)$, $t = 1, 2, \ldots$, converges to a (unique) fixed point $x^*$ satisfying $\mathbf{T}(x^*) = x^*$ starting from any initial value $x^0 \in \mathcal{X}$. To prove Proposition 5, we shall show that the function $h_2(\cdot)$, defined in (3.7), is a contraction mapping on $[0, D_{\max}]$ with respect to the absolute value norm if the condition (3.8) is satisfied.

Let $D_{2,a}$ and $D_{2,b}$ be two arbitrarily-chosen real numbers such that $D_{\max} \geq D_{2,a} > D_{2,b} \geq 0$. Then, it can be shown that

$$\left| h_2(D_{2,a}) - h_2(D_{2,b}) \right| = \int_{y=0}^{d_{\max}} y \left\{ \int_{x=g_2(D_{2,b})+p_2+\gamma_2[y-d_2^*]^+}^{g_2(D_{2,a})+p_2+\gamma_2[y-d_2^*]^+} f(x,y)dx \right\} dy. \tag{3.10}$$

Denote $F(\theta \mid y) = \int_{x=-\infty}^{\theta} f(x,y)dx$. Thus, we can obtain the following inequalities

$$\int_{x=g_2(D_{2,b})+p_2+\gamma_2[y-d_2^*]^+}^{g_2(D_{2,a})+p_2+\gamma_2[y-d_2^*]^+} f(x,y)dx \tag{3.11}$$

$$= F(g_2(D_{2,a}) + p_2 + \gamma_2[y-d_2^*]^+ \mid y) - F(g_2(D_{2,b}) + p_2 + \gamma_2[y-d_2^*]^+ \mid y) \tag{3.12}$$

$$= f(D_{2,\gamma}, y) \cdot g_2'(D_{2,\gamma}) \cdot (D_{2,a} - D_{2,b}) \tag{3.13}$$

$$\leq K \cdot g_2'(D_{2,\gamma}) \cdot (D_{2,a} - D_{2,b}) \tag{3.14}$$

$$\leq K \cdot \left[ \max_{D_2 \in [0,D_{\max}]} g_2'(D_2) \right] \cdot (D_{2,a} - D_{2,b}), \tag{3.15}$$

where (3.13) follows from the intermediate value theorem and chain rule, $D_{2,\gamma}$ is a certain value in $[D_{2,b}, D_{2,a}]$, and $K = \max_{(\theta,d) \in \mathcal{U}} f(\theta, d)$. Then, by plugging the inequality (3.15) into (3.10), we have

$$|h_2(D_{2,a}) - h_2(D_{2,b})| \leq \int_{y=0}^{d_{\max}} yK \left[ \max_{D_2 \in [0,D_{\max}]} g_2'(D_2) \right] (D_{2,a} - D_{2,b})dy \tag{3.16}$$

$$= \frac{d_{max}^2}{2} K \left[ \max_{D_2 \in [0,D_{\max}]} g'(D_2) \right] \cdot |D_{2,a} - D_{2,b}|. \tag{3.17}$$

Therefore, if the condition (3.8) is satisfied, then $\kappa = \frac{d_{max}^2}{2} K \left[ \max_{D_2 \in [0,D_{\max}]} g_2'(D_2) \right] \in [0,1)$ and $|h_2(D_{2,a}) - h_2(D_{2,b})| \leq \kappa |D_{2,a} - D_{2,b}|$, for any $D_{\max} \geq D_{2,a} > D_{2,b} \geq 0$. In other words, $h_2(\cdot)$ is a contraction mapping on $[0, D_{\max}]$ with respect to the absolute value norm. Thus, by applying Proposition 1.1 in Chapter 3 of [57], we see that the data demand dynamics converges if the condition (3.8) is satisfied. Proposition 5 is therefore proved. ∎

Proposition 5 states the relation between the congestion cost function and the distribution of $(\theta, d)$ such that the data demand dynamics converges, and holds for a general yet practical data plan. Although the convergence condition (3.8) is sufficient but not necessary, it provides us with the insight that, for a given distribution function $f(\theta, d)$, if the congestion cost increases too fast (i.e., $g_2'(D_2)$ is larger than $2/(K \cdot d_{\max}^2)$ for some $D_2 \in [0, D_{\max}]$), the data demand dynamics may oscillate or diverge. A simi-

lar insight was reported in [49] in the context of the decentralized spectrum access in cognitive networks. Next, by considering $g_2(D_2) = D_2/C_2$ for $D_2 \in [0, D_{\max}]$, we investigate the impacts of the WSP's long-term capacity decision on the convergence of the users' data demand dynamics. The result is summarized as follows.

**Corollary 4.** *Suppose that $g_2(D_2) = D_2/C_2$ for $D_2 \in [0, D_{\max}]$. For any data plan $\mathcal{P}_2 = (p_2, \gamma_2, d_2^*)$, the data demand dynamics specified by (3.7) converges to the unique equilibrium point starting from any initial point $D_2^0 \in [0, D_{\max}]$ if*

$$C_2 > \frac{K \cdot d_{\max}^2}{2}, \tag{3.18}$$

*where $K = \max_{(\theta,d) \in \mathcal{U}} f(\theta, d)$. If $f(\theta, d) = 1$ over $\mathcal{U} = \{(\theta, d) \,|\, 0 \le \theta \le 1, 0 \le d \le 1\}$, then (3.18) becomes $C_2 > 1/2$.* ∎

Corollary 4 indicates that the network capacity allocated to $\mathcal{P}_2-$users needs to be greater than a certain threshold such that the data demand dynamics is guaranteed to converge for any data plan $\mathcal{P}_2 = (p_2, \gamma_2, d_2^*)$. In particular, if $(\theta, d)$ is uniformly distributed over $\mathcal{U} = \{(\theta, d) \,|\, 0 \le \theta \le 1, 0 \le d \le 1\}$, then the capacity threshold (normalized with respect to the total number of users) corresponds to data demand averaged over all the users in the market. This implies that the network for the plan $\mathcal{P}_2$ needs to be able to accommodate all the users' data demand. Moreover, we see from (3.18) that the capacity threshold $K \cdot d_{\max}^2/2$ does not explicitly depend on $D_{\max}$. Instead, it is closely related to $d_{\max}$. In particular, if $d_{\max}$ increases, then a more stringent requirement is imposed on the WSP's network capacity in order to guarantee the convergence of data demand dynamics regardless of the initial points or data plans. On the other hand, if the network capacity is not large enough (e.g., $C_2 < D_{\max}$), then the users may experience excessive delays (i.e., high congestion costs) and the data demand dynamics may oscillate without convergence.

Before studying the WSP's data plan decision, we make two remarks regarding the

users' subscription decisions.

*Remark 5:* The dynamics specified by (3.7) requires that all the users update subscription decisions at the beginning of each time period. In practice, if only a fraction $\epsilon \in (0, 1]$ of the user update subscription decisions each time, then the sequence becomes

$$D_2^t = \epsilon h_2(D_2^{t-1}) + (1 - \epsilon)D_2^{t-1}, \tag{3.19}$$

where $h_2(D_2^{t-1})$ is given by (3.7). The equilibrium analysis is not affected, whereas the convergence condition in (3.8) is modified as $\max_{D_2 \in [0, D_{\max}]} g_2'(D_2) < \frac{2}{\epsilon \cdot (K \cdot d_{\max}^2)}$, which is more easily satisfied for a smaller $\epsilon \in (0, 1]$. In other words, the parameter $\epsilon \in (0, 1]$ smooths the data demand update process and makes the dynamics easier to converge by slowing down the convergence rate. Nevertheless, if the network capacity is large enough to serve all the users' data demand in practice, then the convergence can always be observed even though all the users update their subscription decisions. It should also be noted that another approach to modeling the users' data demand dynamics is considering a continuous-time dynamics specified as

$$\frac{dD_2}{dt} = \rho \cdot [h_2(D_2) - D_2], \tag{3.20}$$

where $\rho$ is referred to as the diffusion rate [66]. For (3.20), the equilibrium is still defined the same as that in Definition 1, while the convergence is guaranteed. The considered discrete-time data demand dynamics has been studied in prior work (see, e.g., [38][42]) and is more appropriate for scenarios in which the users' subscription decisions can only change in discrete time instants (e.g., at the beginning of a day or month). Moreover, $\epsilon$ in (3.19) is essentially the same as $\rho$ in (3.20) and the discrete-time dynamics considered in this chapter will become (3.20) if the duration of a time period is sufficiently small.

*Remark 6:* As in the existing literature [35]–[39], the cost in updating the subscription decisions (e.g., time spent in calling the customer service, activation fees and early termination fees) are not considered in the chapter. Here, we briefly discuss the impacts of this cost on the data demand dynamics. For simplicity, we assume that the cost of activating the data plan and that of terminating the subscription are the same, and we refer to this cost as *switching* cost denoted by $c_s$. With a switching cost, the users' subscription decisions are affected. Specifically, if user $k$ is a subscriber in the time period $t$, it will continue the subscription in the next time period $t+1$ if

$$\theta_k - g_2(D_2^t) - p_2 - \gamma_2[d_k - d_2^*]^+ \geq -c_s. \tag{3.21}$$

On the other hand, if user $k$ is not a subscriber in the time period $t$, it will choose to subscribe to the data plan in the next time period $t+1$ if

$$\theta_k - g_2(D_2^t) - p_2 - \gamma_2[d_k - d_2^*]^+ - c_s \geq 0. \tag{3.22}$$

It should be noted that if the cost is taken into account when the users make their subscription decisions, there may exist multiple equilibrium data demand points, and the convergence is subject to the initial point. For instance, in the extreme case in which the cost is so high (e.g., greater than $\theta_{\max}$) that no users would like to update their subscription decisions, every possible value of (aggregate) data demand $D_2 \in [0, D_{\max}]$ is an equilibrium point. We shall show in the numerical results the impact of switching cost $c_s$ on the users' subscription decisions, while rigorous analysis of $c_s$ is left as our future work.

### 3.4.2 WSP's Data Plan Decision

Over the entire lifespan of the network infrastructure, the WSP can change its data plans to maximize its revenue, although the change of data plans is sufficiently slow compared to the users' subscription decisions. In other words, the duration of a medium-term period corresponds to that of a sufficiently large number of short-term periods. We note that given the WSP's data plan, the data demand dynamics converges rapidly (e.g., within a few iterations) to the equilibrium point if the convergence condition (3.8) is satisfied. Thus, the WSP's average revenue per short-term period (i.e., users' subscription period) is approximately equal to its equilibrium revenue when the data demand reaches the unique equilibrium. Next, we derive the expression of the WSP's equilibrium revenue as follows.

$$
\begin{aligned}
R_2 =& p_2 \int_{y=0}^{d_{\max}} \int_{x=g_2(D_2^*)+p_2+\gamma_2[y-d_2^*]^+}^{\theta_{\max}} f(x,y)dxdy \\
&+ \int_{y=d_2^*}^{d_{\max}} \int_{x=g_2(D_2^*)+p_2+\gamma_2(y-d_2^*)}^{\theta_{\max}} \gamma_2 \cdot (y - d_2^*) \cdot f(x,y)dxdy,
\end{aligned}
\tag{3.23}
$$

where the first term on the right hand side is the subscription fee that every subscriber pays and the second term is the additional fee that users with demand higher than $d_2^*$ pays. Although the equilibrium data demand $D_2^*$ is uniquely determined by the WSP's data plan $\mathcal{P}_2 = (p_2, d_2^*, \gamma_2)$ and hence can be expressed as an implicit function of $\mathcal{P}_2$, it is rather challenging to maximize the equilibrium revenue in (3.23). The difficulties are mainly: (1) $D_2^*$ cannot be expressed explicitly in a closed-form function in terms of $\mathcal{P}_2 = (p_2, d_2^*, \gamma_2)$; (2) due to the integral, the equilibrium revenue is not an explicit function of $\mathcal{P}_2 = (p_2, \gamma_2, d_2^*)$. Thus, we resort to numerical methods to find the optimal $\mathcal{P}_2 = (p_2, d_2^*, \gamma_2)$ maximizing the equilibrium revenue. Specifically, we search over all the possible values of $(p_2, d_2^*, \gamma_2)$ and select the one that yields the maximum equilibrium revenue. In practice, the data plan is typically confined within a small

66

finite set of options[9] and hence, the complexity associated with the exhaustive search is not prohibitive.

In the following, we propose a heuristic algorithm that progressively chooses the (locally) optimal data plan in a greedy manner. For the ease of presenting the algorithm and deriving more specific results, we consider uniformly distributed $(\theta, d)$, i.e., $f(\theta, d) = 1$ over $\mathcal{U} = \{(\theta, d) \mid 0 \leq \theta \leq 1, 0 \leq d \leq 1\}$, although other forms of $f(\theta, d)$ can also be applied. Under the assumption of uniformly distributed of $(\theta, d)$ over $\mathcal{U} = \{(\theta, d) \mid 0 \leq \theta \leq 1, 0 \leq d \leq 1\}$, we rewrite the equilibrium revenue in (3.23) as

$$R_2 = p_2 \left[1 - g_2(D_2^*) - p_2\right] + \frac{\gamma_2(1 - d_2^*)^2}{2} \left[1 - g_2(D_2^*) - 2p_2 - \frac{2\gamma_2(1 - d_2^*)}{3}\right]$$

(3.24)

where $p_2 \leq 1 - g_2(D_2^*)$, $d_2^* \leq 1$ and $\gamma_2 \geq 0$.[10] Note that, even if we artificially assume that the congestion cost $g_2(D_2^*)$ is independent of the data plan, (3.24) is non-concave in $(p_2, \gamma_2, d_2^*)$. Thus, there exist no efficient algorithms to find the optimal $(p_2, d_2^*, \gamma_2)$. In the proposed heuristic algorithm, instead of jointly optimizing $(p_2, d_2^*, \gamma_2)$, we optimize $p_2$, $d_2^*$, and $\gamma_2$ separately. Specifically, by assuming that the equilibrium data demand $D_2^*$ is independent of $p_2$ and treating $g_2(D_2^*)$, $d_2^*$ and $\gamma_2$ as fixed values, we choose the optimal $p_2$ to maximize (3.24). Then, we apply the same technique to optimize $d_2^*$ and $\gamma_2$, and the same process repeats until the stopping criterion is satisfied (e.g., convergence or the maximum number of iterations is reached). To summarize, the heuristic algorithm is described in Algorithm I.

---

[9]In practice, the subscription fee $p_2$ is usually selected from $\{9.99\$, 19.99\$, 24.99\$, 29.99\$\}$ or similar set of options.

[10]Since the maximum demand is $d_{\max} = 1$, $d_2^* > 1$ and $d_2^* = 1$ are essentially the same.

**Algorithm 1** Find $(p_2, d_2^*, \gamma_2)$

---

$R_2 = 0$, $d_2^* \leftarrow \infty$, $\gamma_2 \leftarrow 0$, and $k \leftarrow 1$
**while** $k \leq MaxIterate$ **do**
  $temp \leftarrow R_2$
  Optimize $p_2$: $p_2 \leftarrow \left[ \frac{1 - g_2(D_2^*) - \gamma_2(1 - d_2^*)^2}{2} \right]^+$
  Optimize $d_2^*$: $d_2^* \leftarrow \left[ \frac{1 - g_2(D_2^*) - p_2}{\gamma_2} \right]_0^1$
  Optimize $\gamma_2$: $\gamma_2 \leftarrow \left[ \frac{3 \left[ 1 - g_2(D_2^*) - p_2 \right]}{4(1 - d_2^*)} \right]^+$
  Recalculate $R_2$ based on (3.24)
  **if** abs$(temp - R_2) \leq$ threhold **then**
    break
  **end if**
  Update $D_2^*$, $g_2(D_2^*)$ and $k++$
**end while**
**return** $(p_2, d_2^*, \gamma_2)$

---

### 3.4.3 WSP's Capacity Decision

We assume that the WSP's network capacity is chosen to guarantee the convergence of data demand dynamics to the unique equilibrium point regardless of the initial points.[11] For instance, if $(\theta, d)$, i.e., $f(\theta, d) = 1$ is uniformly distributed over $\mathcal{U} = \{(\theta, d) \,|\, 0 \leq \theta \leq 1, 0 \leq d \leq 1\}$ and $g_2(D_2) = D_2/C_2$, then the network capacity allocated to $\mathcal{P}_2-$users should be greater than $1/(2\epsilon)$, where $\epsilon$ is the fraction of users that update their subscription decisions in each time period. As can be seen from (3.23), given the WSP's capacity, it is rather difficult to find explicitly the optimal value of $(p_2, d_2^*, \gamma_2)$ maximizing the WSP's equilibrium revenue. As a result, we can only numerically find the optimal network capacity to maximize the WSP's equilibrium profit in (3.1).

Finally, we note that if only the unlimited data plan $\mathcal{P}_1 = \{p_1, +\infty, 0\}$ is offered, the above analysis still applies and the corresponding result can be easily obtained by letting $d_2^* = \infty$ and $\gamma_2 = 0$.

---

[11]This requires that the average cost of $\tau$ per unit capacity be sufficiently small such that the WSP can receive a non-negative profit.

## 3.5 Wireless Communications Market: Two Data Plans

In this part, we turn to the analysis of a wireless communications market where the WSP offers two data plans $\mathcal{P}_1$ and $\mathcal{P}_2$. Although we mainly focus on the scenario that these two data plans are offered by the same monopolistic WSP, we shall also briefly discuss the case in which they are offered by two competing WSPs.

### 3.5.1 Users' Subscription Decisions

As in a market with only one data plan $\mathcal{P}_2$, we study the users' equilibrium subscription decisions by specifying the equilibrium data demand $(D_1^*, D_2^*)$. By Assumption 2, we see that the equilibrium data demand $(D_1^*, D_2^*)$ satisfies the following equations

$$D_1^* = h_{d,1}(D_1^*, D_2^*) = \int_{y=\tilde{d}}^{d_{\max}} \int_{x=g_1(D_1^*)+p_1}^{\theta_{\max}} y f(x,y) dx dy \tag{3.25}$$

$$D_2^* = h_{d,2}(D_1^*, D_2^*) = \int_{y=0}^{\tilde{d}} \int_{x=g_2(D_2^*)+p_2+\gamma_2[y-d_2^*]^+}^{\theta_{\max}} y f(x,y) dx dy \tag{3.26}$$

if $p_1 + g_1(D_1^*) > p_2 + g_2(D_2^*)$, and

$$D_1^* = h_{d,1}(D_1^*, D_2^*) = \int_{y=0}^{d_{\max}} \int_{x=g_1(D_1^*)+p_1}^{\theta_{\max}} y f(x,y) dx dy \tag{3.27}$$

$$D_2^* = h_{d,2}(D_1^*, D_2^*) = 0 \tag{3.28}$$

if $p_1 + g_1(D_1^*) \leq p_2 + g_2(D_2^*)$. In (3.25) and (3.26), $\tilde{d}$ is given by

$$\tilde{d} = d_2^* + \frac{1}{\gamma_2}[p_1 - p_2 + g_1(D_1^*) - g_2(D_2^*)], \tag{3.29}$$

which specifies the data demand of marginal users that are "indifferent" between subscribing to the plan $\mathcal{P}_1$ and the plan $\mathcal{P}_2$ (see [35][39] for a detailed explanation of

"indifferent"). Note that there are two regimes of the equilibrium data demand in the market with two data plans, and which regime governs the equilibrium depends on the relative values of the *effective* full price (not including the additional cost if the data demand exceeds the granted data limit), i.e., $p_1 + g_1(D_1^*)$ and $p_2 + g_2(D_2^*)$. Next, we give the formal definition of the equilibrium point $(D_1^*, D_2^*)$, which is similar to Definition 1.

*Definition 2:* When two data plans $\mathcal{P}_1$ and $\mathcal{P}_2$ are offered, $(D_1^*, D_2^*)$ is an *equilibrium* data demand if it satisfies satisfies

$$h_{d,1}(D_1^*, D_2^*) = D_1^* \text{ and } h_{d,2}(D_1^*, D_2^*) = D_2^*, \tag{3.30}$$

where $h_{d,1}(D_1^*, D_2^*)$ and $h_{d,2}(D_1^*, D_2^*)$ are given in (3.25)–(3.28).

When the unlimited data plan $\mathcal{P}_1$ is available in the market, there may not exist an equilibrium data demand if the plan $\mathcal{P}_2$ is "capped" (i.e., $d_2^* < d_{\max}$ and $\gamma_2 > 0$). Suppose, as an counter example, that $g_2(D_2) = 0$ is a constant for $D_2 \in [0, D_{\max}]$. Thus, $h_{d,1}(D_1, D_2)$ in (3.25) and (3.27) is independent of $D_2$ and can be rewritten compactly as $h_{d,1}(D_1)$. From (3.25) and (3.27), we see that the integration interval is not continuous, implying that $\tilde{h}_{d,1}(D_1) = h_{d,1}(D_1) - D_1$ may not be a continuous function in $D_1 \in [0, D_{\max}]$. Specifically, if $p_1 + g_1(D_1) \leq p_2$, then the integration interval is $[0, d_{\max}]$, i.e., no users subscribes to the plan $\mathcal{P}_2$, whereas if $p_1 + g_1(D_1) > p_2$, the integration interval is $[\tilde{d}, d_{\max}]$. According to Definition 2, the equilibrium data demand should satisfy $\tilde{h}_{d,1}(D_1^*) = h_{d,1}(D_1^*) - D_1^* = 0$. Although it is easy to show that $\tilde{h}_{d,1}(D_1)$ is strictly decreasing in $D_1 \in [0, D_{\max}]$, $\tilde{h}_{d,1}(0) \geq 0$ and $\tilde{h}_{d,1}(D_{\max}) \leq 0$, it is not guaranteed that $\tilde{h}_{d,1}(D_1)$ has a root, since $\tilde{h}_{d,1}(D_1)$ may not be a continuous function in $D_1 \in [0, D_{\max}]$. In other words, an equilibrium data demand may not exist. Next, we provide a sufficient condition that establishes the existence and uniqueness of an equilibrium point in Proposition 6.

**Proposition 6.** *For any data plans $\mathcal{P}_1 = (p_1, +\infty, 0)$ and $\mathcal{P}_2 = (p_2, d_2^*, \gamma_2)$, there exists a unique equilibrium data demand $(D_1^*, D_2^*)$ satisfying (3.25)–(3.28) if*

$$d_2^* = 0 \text{ and } \gamma_2 > 0. \tag{3.31}$$

*Moreover, the equilibrium data demand $(D_1^*, D_2^*)$ satisfies $D_1^* = h_{d,1}(D_1, 0^*)$ and $D_2^* = 0$ if $p_2 + g_2(0) \geq p_1 + g_1(D_1^*)$.*

*Proof.* To facilitate the proof, we first define two auxiliary functions $\tilde{h}_{d,1}(D_1, D_2) = h_{d,1}(D_1, D_2) - D_1$ and $\tilde{h}_{d,2}(D_1, D_2) = h_{d,2}(D_1, D_2) - D_2$ for $(D_1, D_2) \in \mathcal{D}$, where $h_{d,1}(D_1, D_2)$ and $h_{d,2}(D_1, D_2)$ are defined in (3.25)–(3.28). By Definition 2, $(D_1^*, D_2^*)$ is an equilibrium point if and only if

$$\tilde{h}_{d,1}(D_1^*, D_2^*) = 0 \text{ and } \tilde{h}_{d,2}(D_1^*, D_2^*) = 0. \tag{3.32}$$

Hence, it suffices to show that the equation set in (3.32) has a unique solution on its domain $\mathcal{D}$.

Let us first assume that $D_2 \in [0, D_{\max}]$ is a fixed number. We can show that, if $\gamma_2 > 0$ and $d_2^* = 0$, $\tilde{h}_{d,1}(D_1, D_2)$ is a strictly decreasing and continuous function of $D_1 \in [0, D_{\max} - D_2]$. Moreover, for any fixed value of $D_2 \in [0, D_{max}]$, we have $\tilde{h}_{d,1}(D_{\max} - D_2, D_2) \leq 0$ and $\tilde{h}_{d,1}(0, D_2) \geq 0$. Therefore, by applying the intermediate value theorem, it follows that $\tilde{h}_{d,1}(D_{\max} - D_2, D_2)$ has a unique root $D_1^*$ given any fixed value of $D_2$. Thus, $D_1^*$ can be expressed as a function in terms of $D_2$, and $\tilde{h}_{d,1}(D_1^*, D_2)$ and $\tilde{h}_{d,2}(D_1^*, D_2)$ can be rewritten in a compact form as $\tilde{h}_{d,1}(D_2)$ and $\tilde{h}_{d,2}(D_2)$, respectively. It can also be easily proved that $D_{1*}$ is a decreasing function of $D_2$ in $D_2 \in [0, D_{\max}]$. Next, we need to show that $\tilde{h}_{d,2}(D_2)$ has a unique root in $D_2 \in [0, D_{\max}]$ in order to prove Proposition 6.

**Lemma 1.** $h_t(D_2) = h_{d,2}(D_1^*, D_2) + D_1^*$ *is decreasing in $D_2$.*

*Proof.* Note that $h_{d,2}(D_1^*, D_2) + D_1^*$ is the sum data demand of WSP $W_1$ operating at its equilibrium data demand point and $W_2$ at the next period, when WSP $W_2$ currently has a data demand of $D_2$. Thus, $h_{d,2}(D_1^*, D_2) + D_1^*$ can be expressed as follows

$$
\begin{aligned}
h_{d,2}(D_1^*(D_2), D_2) + D_1^*(D_2) = D_{max} &- \int_{y=0}^{\tilde{d}} \int_{x=0}^{g_2(D_2)+p_2+\gamma_2 y} y f(x,y) dx dy \\
&- \int_{y=\tilde{d}}^{d_{\max}} \int_{x=0}^{g_1(D_1^*)+p_1} y f(x,y) dx dy,
\end{aligned}
\tag{3.33}
$$

where $\tilde{d} = \frac{1}{\gamma_2} [p_1 - p_2 + g_1(D_1^*) - g_2(D_2)]^+$, the second term and third term on the right side of the equality represent the aggregate data demand of those users that do not subscribe to either WSP. Since $D_1^*(D_2)$ is increasing in $D_2 \in [0, D_{\max}]$ and $g_1(D_1)$ is increasing in $D_1$ in its domain, we see that $g_1(D_1^*) = g_1(D_1^*(D_2))$ is also increasing in $D_2 \in [0, D_{\max}]$. Therefore, (3.33) is decreasing in $D_2 \in [0, D_{\max}]$. This can also be intuitively expected. When the data demand of both WSPs increases, the congestion costs increase and hence fewer users will subscribe to the WSPs, which will in turn result in a decrease in the total data demand of these two WSPs. $\qquad\square$

Recall that $D_1^*(D_2)$ is increasing in $D_2 \in [0, D_{\max}]$. Thus, following Lemma 1, it can be seen that $h_{d,2}(D_1^*(D_2), D_2) = h_t(D_2) - D_1^*(D_2)$ is a non-increasing function of $D_2$ and $\tilde{h}_{d,2}(D_1^*, D_2) = \tilde{h}_{d,2}(D_2)$ is a strictly decreasing function of $D_2$ in $[0, D_{\max}]$. On the one hand, $h_{d,2}(D_{\max}) - D_{\max} \le 0$ and, on the other hand, $h_{d,2}(0) - 0 \ge 0$. Thus, due to its continuity and strictly decreasing property, $\tilde{h}_{d,2}(D_1^*, D_2) = h_{d,2}(D_1^*(D_2), D_2) - D_2$ has a unique root in $D_2 \in [0, D_{\max}]$. This proves Proposition 6. $\qquad\blacksquare$

Proposition 6 indicates that, if the two data plans $\mathcal{P}_1$ and $\mathcal{P}_2$ are unlimited and usage-based, respectively, then the data demand admits a unique equilibrium point. It also shows that, if the effective subscription cost of for the data plan $\mathcal{P}_1$ evaluated at $D_1^*$ is always smaller than or equal to that of the data plan $\mathcal{P}_2$, then no users subscribe to the data plan $\mathcal{P}_2$ at the equilibrium point.

Following Chapter 3.4, we consider a discrete-time best-response dynamics to model the users' subscription decision process. With two data plans $\mathcal{P}_1 = (p_1, +\infty, 0)$ and $\mathcal{P}_2 = (p_2, d_2^*, \gamma_2)$ offered in the market, each user has three possible choices at the beginning of each time period: subscribe to the plan $\mathcal{P}_1$, subscribe to the plan $\mathcal{P}_2$, and subscribe to neither. The users expect that the congestion cost incurred when subscribing to a data plan in the time period $t$ is equal to that in the previous period $t - 1$ and make their subscription decisions to myopically maximize their utility in the time period $t$ [35][66][38]. We assume that, other than the subscription price, there is no cost involved (e.g., initiation fees, termination fees, device prices) when users switch between the data plans $\mathcal{P}_1$ and $\mathcal{P}_2$[66]. By Assumption 2, at period $t = 1, 2 \cdots$, user $k$ subscribes to the data plan $\mathcal{P}_1$ if and only if

$$\theta_k - g_1(D_1^{t-1}) - p_1 \geq \theta_k - g_2(D_2^{t-1}) - p_2 - \gamma_2[d_k - d_2^*]^+ \text{ and } \theta_k - g_1(D_1^{t-1}) - p_1 \geq 0,$$

(3.34)

to the data plan $\mathcal{P}_2$ if and only if

$$\theta_k - g_2(D_2^{t-1}) - p_2 - \gamma_2[d_k - d_2^*]^+ > \theta_k - g_1(D_1^{t-1}) - p_1$$
$$\text{and } \theta_k - g_2(D_2^{t-1}) - p_2 - \gamma_2[d_k - d_2^*]^+ \geq 0,$$

(3.35)

and to neither data plan if and only if

$$\theta_k - g_1(D_1^{t-1}) - p_1 < 0 \text{ and } \theta_k - g_2(D_2^{t-1}) - p_2 - \gamma_2[d_k - d_2^*]^+ < 0. \quad (3.36)$$

Therefore, given the data plans $\mathcal{P}_1 = (p_1, +\infty, 0)$ and $\mathcal{P}_2 = (p_2, \gamma_2, d_2^*)$, the data demand dynamics is described by a sequence $\{(D_1^t, D_2^t)\}_{t=0}^{\infty}$ in $\mathcal{D} = \{(D_1, D_2) \in \mathbb{R}_+^2 \mid D_1 + D_2 \leq D_{\max}\}$ generated by $D_1^t = h_{d_1}(D_1^{t-1}, D_2^{t-1})$ and $D_2^t = h_{d_2}(D_1^{t-1}, D_2^{t-1})$, where $h_{d_1}(D_1^{t-1}, D_2^{t-1})$ and $h_{d_2}(D_1^{t-1}, D_2^{t-1})$ are obtained by substituting $(D_1^{t-1}, D_2^{t-1})$ into (3.25)–(3.28).

Since an equilibrium point may not exist if the data plan $\mathcal{P}_2$ is unlimited or capped, we restrict the analysis in the remainder of this chapter to the case that the plan $\mathcal{P}_2$ is usage-based (although an initial subscription fee $p_2$ may be charged) such that a unique equilibrium point is guaranteed to exist. Next, we provide a sufficient condition for the data demand dynamics $\{(D_1^t, D_2^t)\}_{t=0}^{\infty}$ to converge.

**Proposition 7.** *For data plans $\mathcal{P}_1 = (p_1, +\infty, 0)$ and $\mathcal{P}_2 = (p_2, 0, \gamma_2)$ where $\gamma_2 > 0$, the data demand dynamics converges to the unique equilibrium point starting from any initial point $(D_1^0, D_2^0) \in \mathcal{D} = \{(D_1, D_2) \in \mathbb{R}_+^2 \mid D_1 + D_2 \leq D_{\max}\}$ if the following condition is satisfied:*

$$\max_{(D_1, D_2) \in [0, d_{\max}]^2} \{g_1'(D_1), g_2'(D_2)\} < \frac{1}{K \cdot \left( \frac{d_{\max}^2}{2} + \frac{d_{\max}}{\gamma_2}[\theta_{\max} - p_1]^+ + \frac{d_{\max}}{\gamma_2}[\theta_{\max} - p_2]^+ \right)},$$

(3.37)

*where $d_{max}$ is the maximum individual demand, $\theta_{\max}$ is the maximum benefit derived from subscribing to the WSP's service and $K = \max_{(\theta, d) \in \mathcal{U}} f(\theta, d)$*

*Proof.* First, define the mapping that specifies the data demand dynamics by $\mathbf{h}_d$ :

$$\mathbf{h}_d(D_1, D_2) = \Big( h_{d,1}(D_1, D_2), h_{d,2}(D_1, D_2) \Big),$$

(3.38)

where $h_{d,1}$ and $h_{d,2}$ are defined in (3.25)–(3.28). In order to establish the global convergence of the data demand dynamics, we shall show that the mapping $\mathbf{h}_d(\cdot)$ is a contraction on $\mathcal{D}$ with respect to a certain norm [57]. Unlike in a market with only one data plan offered, the mapping $\mathbf{h}_d$ is no longer a scalar function and hence the absolute value norm is not applicable. Instead, we apply $L_1$ norm and show that $\mathbf{h}_d(\cdot)$ is a contraction with respect to $L_1$ norm if the condition (3.37) is satisfied.

Let $D_a = (D_{1,a}, D_{2,a}) \in \mathcal{D}$ and $D_b = (D_{1,b}, D_{2,b}) \in \mathcal{D}$ be two arbitrarily-chosen

points such that $\tilde{d}_a \geq \tilde{d}_b \geq 0$, where $\tilde{d}_a = \frac{1}{\gamma_2} \left[ p_1 - p_2 + g_1(D_{1,a}) - g_2(D_{2,a}) \right]^+$ and $\tilde{d}_b = \frac{1}{\gamma_2} \left[ p_1 - p_2 + g_1(D_{1,b}) - g_2(D_{2,b}) \right]^+$. By the definition of $L_1$ norm, we have

$$
\left\| \mathbf{h}_d(D_{1,a}, D_{2,a}) - \mathbf{h}_d(D_{1,b}, D_{2,b}) \right\|_1
$$
$$
= \left| h_{d,1}(D_{1,a}, D_{2,a}) - h_{d,1}(D_{1,b}, D_{2,b}) \right| + \left| h_{d,2}(D_{1,a}, D_{2,a}) - h_{d,2}(D_{1,b}, D_{2,b}) \right|.
\tag{3.39}
$$

Note that the term $\left| h_{d,1}(D_{1,a}, D_{2,a}) - h_{d,1}(D_{1,b}, D_{2,b}) \right|$ can be expanded and rewritten as

$$
\left| h_{d,1}(D_{1,a}, D_{2,a}) - h_{d,1}(D_{1,b}, D_{2,b}) \right|
\tag{3.40}
$$
$$
= \left| \int_{y=\tilde{d}_a}^{d_{\max}} \int_{x=g_1(D_{1,a})+p_1}^{\theta_{\max}} yf(x,y)dxdy - \int_{y=\tilde{d}_b}^{\tilde{d}_a} \int_{x=g_1(D_{1,b})+p_1}^{\theta_{\max}} yf(x,y)dxdy \right.
$$
$$
\left. - \int_{y=\tilde{d}_a}^{d_{\max}} \int_{x=g_1(D_{1,b})+p_1}^{\theta_{\max}} yf(x,y)dxdy \right|
\tag{3.41}
$$
$$
= \left| \int_{y=\tilde{d}_a}^{d_{\max}} \int_{x=g_1(D_{1,a})+p_1}^{g_1(D_{1,b})+p_1} yf(x,y)dxdy - \int_{y=\tilde{d}_b}^{\tilde{d}_a} \int_{x=g_1(D_{1,b})+p_1}^{\theta_{\max}} yf(x,y)dxdy \right|
\tag{3.42}
$$

Denote $K = \max_{(\theta,d)\in\mathcal{U}} f(\theta, d)$. Next, we show that the following inequalities can be established

$$
\left| \int_{y=\tilde{d}_a}^{d_{\max}} \int_{x=g_1(D_{1,a})+p_1}^{g_1(D_{1,b})+p_1} yf(x,y)dxdy - \int_{y=\tilde{d}_b}^{\tilde{d}_a} \int_{x=g_1(D_{1,b})+p_1}^{\theta_{\max}} yf(x,y)dxdy \right|
$$
$$
\leq K \cdot \left| g_1(D_{1,a}) - g_1(D_{1,b}) \right| \frac{d_{\max}^2 - \min\{d_{\max}^2, \tilde{d}_a^2\}}{2} + K \cdot \left[ \theta_{\max} - (p_1 + g_1(D_{1,b})) \right]^+ \frac{\tilde{d}_a^2 - \tilde{d}_b^2}{2}
\tag{3.43}
$$
$$
\leq K \cdot g_1'(D_{1,c}) \cdot \left| D_{1,a} - D_{1,b} \right| \cdot \frac{d_{\max}^2}{2} + K \cdot \left[ \theta_{\max} - p_1 \right]^+ \cdot \frac{(\tilde{d}_a + \tilde{d}_b)(\tilde{d}_a - \tilde{d}_b)}{2}
\tag{3.44}
$$
$$
\leq K \cdot g_1'(D_{1,c}) \cdot \left| D_{1,a} - D_{1,b} \right| \cdot \frac{d_{\max}^2}{2} + K \cdot \left[ \theta_{\max} - p_1 \right]^+ \cdot d_{\max} \cdot \left[ \tilde{d}_a - \tilde{d}_b \right]^+
\tag{3.45}
$$
$$
\leq K \cdot g_1'(D_{1,c}) \cdot \left| D_{1,a} - D_{1,b} \right| \cdot \frac{d_{\max}^2}{2} + K \cdot \left[ \theta_{\max} - p_1 \right]^+ \cdot \frac{d_{\max}}{\gamma_2} \cdot \left| g_1(D_{1,a}) - g_1(D_{1,b}) \right|
$$
$$
- K \cdot \left[ \theta_{\max} - p_1 \right]^+ \cdot \frac{d_{\max}}{\gamma_2} \left| g_2(D_{2,a}) - g_2(D_{2,b}) \right|,
\tag{3.46}
$$

where $D_{1,c}$ is a number between $D_{1,a}$ and $D_{1,b}$, (3.44) follows (3.43) based on the intermediate value theorem, and (3.46) is due to the fact that $\left[[x_1]^+ - [x_2]^+\right]^+ \leq [x_1 - x_2]^+$. Thus, by combining (3.40)–(3.46), we see that (3.40) is less than or equal to (3.46).

Similarly, it can be shown that

$$\left| h_{d,2}(D_{1,a}, D_{2,a}) - h_{d,2}(D_{1,b}, D_{2,b}) \right| \tag{3.47}$$

$$\leq \quad K \cdot g_2'(D_{2,c'}) \cdot \left| D_{2,a} - D_{2,b} \right| \cdot \frac{d_{\max}^2}{2} + K \cdot \left[ \theta_{\max} - p_2 \right]^+ \cdot \frac{d_{\max}}{\gamma_2} \cdot \left| g_1(D_{1,a}) - g_1(D_{1,b}) \right|$$

$$- K \cdot \left[ \theta_{\max} - p_2 \right]^+ \cdot \frac{d_{\max}}{\gamma_2} \left| g_2(D_{2,a}) - g_2(D_{2,b}) \right|. \tag{3.48}$$

Thus, following (3.40) –(3.48), we have

$$\left\| \mathbf{h}_d(D_{1,a}, D_{2,a}) - \mathbf{h}_d(D_{1,b}, D_{2,b}) \right\|_1 \tag{3.49}$$

$$\leq \quad K \left( \frac{d_{\max}^2}{2} + \frac{d_{\max}}{\gamma_2} \left[ \theta_{\max} - p_1 \right]^+ + \frac{d_{\max}}{\gamma_2} \left[ \theta_{\max} - p_2 \right]^+ \right) \cdot g_1'(D_{1,c}) \cdot \left| D_{1,a} - D_{1,b} \right|$$

$$+ K \left( \frac{d_{\max}^2}{2} + \frac{d_{\max}}{\gamma_2} \left[ \theta_{\max} - p_1 \right]^+ + \frac{d_{\max}}{\gamma_2} \left[ \theta_{\max} - p_2 \right]^+ \right) \cdot g_2'(D_{2,c'}) \cdot \left| D_{2,a} - D_{2,b} \right|.$$

Thus, if the condition in (3.37) is satisfied, the mapping $\mathbf{h}_d(\cdot)$ is a contraction on $\mathcal{D}$ with respect to $L_1$ norm with a modulus $\kappa_d \in [0, 1)$, where $\kappa_d = \max_{(D_1, D_2) \in [0, d_{max}]^2} \left\{ g_1'(D_1), g_2'(D_2) \right\} \cdot K \cdot \left( \frac{d_{\max}^2}{2} + \frac{d_{\max}}{\gamma_2} \left[ \theta_{\max} - p_1 \right]^+ + \frac{d_{\max}}{\gamma_2} \left[ \theta_{\max} - p_2 \right]^+ \right)$, and the data demand dynamics specified by (3.25)–(3.28) converges, regardless of the initial points, to the unique equilibrium point. This proves Proposition 7. ∎

We can obtain more specific condition regarding the network capacities for the convergence of data demand dynamics by plugging $g_1(D_1) = D_1/C_1$ and $g_2(D_2) = D_2/C_2$ into (3.37). The result is similar to Corollary 3.18 and omitted for brevity. Note that the condition (3.37) imposes a more stringent requirement on the congestion costs

(e.g., the WSP needs to allocate larger capacities to the subscribers) than the condition (3.8) does. However, the condition (3.37) provides us with a similar insight that, if congest costs increase too rapidly, the data demand dynamics may exhibit oscillation or divergence. Another important observation from (3.37) is that the two data plans also affect the convergence. Specifically, given higher prices, it is easier for the congestion costs to satisfy the convergence condition. Intuitively, higher prices result in lower aggregate data demand. Therefore, there is less fluctuation in the data demand dynamics and the requirement on the congestion costs becomes less stringent.

### 3.5.2  WSP's Data Plan Decision

Following Chapter 3.4, we first write the the WSP's equilibrium revenues for the data plans $\mathcal{P}_1$ and $\mathcal{P}_2$ as

$$
\begin{aligned}
R_1 &= \int_{y=\tilde{d}}^{d_{\max}} \int_{x=g_1(D_1^*)+p_1}^{\theta_{\max}} p_1 f(x,y) dx dy \\
\text{and } R_2 &= \int_{y=0}^{\tilde{d}} \int_{x=g_2(D_2^*)+p_2+\gamma_2 y}^{\theta_{\max}} (p_2 + \gamma_2 y) f(x,y) dx dy
\end{aligned}
\tag{3.50}
$$

if $p_1 + g_1(D_1^*) > p_2 + g_2(D_2^*)$, and as

$$
R_1 = \int_{y=0}^{d_{\max}} \int_{x=g_1(D_1^*)+p_1}^{\theta_{\max}} p_1 f(x,y) dx dy \text{ and } R_2 = 0
\tag{3.51}
$$

if $p_1 + g_1(D_1^*) \leq p_2 + g_2(D_2^*)$, where $\tilde{d}$ is given by $\tilde{d} = \frac{1}{\gamma_2}[p_1 - p_2 + g_1(D_1^*) - g_2(D_2^*)]$. The expressions of equilibrium revenues in (3.50) and (3.51) are even more complicated than (3.23) and hence, lose analytical tractability. As a consequence, we resort to numerical search to identify the optimal $\mathcal{P}_1 = (p_1, +\infty, 0)$ and $\mathcal{P}_2 = (p_2, 0, \gamma_2)$ maximizing $R_1 + R_2$.

Figure 3.1: Single data plan: oscillation and convergence of data demand dynamics. $\epsilon = 0.5$, $d_2^* = 0.5$, $\gamma_2 = 0.2$. $C_2 = 1/8$ in upper plot and $C_2 = 2/3$ in lower plot.

### 3.5.3 WSP's Capacity Decision

It is mathematically challenging to analytically find the optimal capacities $\mathbf{C} = (C_1, C_2)$ to maximize the WSP's profit, since the optimal data plans can only be numerically found. Thus, as in Chapter 3.4, we find the WSP's optimal capacities through exhaustive search.

In the above analysis, we have considered that the two data plans $\mathcal{P}_1$ and $\mathcal{P}_2$ are offered by the same WSP. Nevertheless, in a wireless communications market, it is possible that these two plans are offered by two different WSPs competing against each other (i.e., duopoly market). The order of moves is almost the same as that described at the beginning of Chapter 3.4, with the exception that in the long-term and medium-term periods, each of the two WSPs decide their own network capacities and data plans, respectively. Specifically, for the long-term capacity decision, the two WSPs simultaneously and independently invest in the network capacities. Then,

Figure 3.2: Single data plan: comparison between discrete-time and continuous-time data demand dynamics. $\rho = 1$, $\epsilon = 0.5$, $d_2^* = 0.5$, $\gamma_2 = 0.2$. $C_2 = 1/8$ in upper plot and $C_2 = 2/3$ in lower plot.

given the capacity decisions, the two WSPs play a noncooperative subgame in which they strategically make data plan decisions. Best-response dynamics can be applied to model the two WSPs' data plan decision process. That is, given its competitor's data plan, each WSP chooses an optimal data plan to selfishly maximize its revenue. In the short-term period, the users' subscription dynamics is unaffected and the same as that studied in Chapter 3.5. Unfortunately, it is mathematically intractable to analyze the competition between the two WSPs, as explicitly expressing the optimal decisions of the two WSPs in response to each other's decision is not possible. With a simpler model, some (partial) analytical results regarding the competition between the WSPs are available in [35][41], whereas in this chapter, we shall illustrate the WSP competition through numerical results.

Figure 3.3: Single data plan: oscillation and convergence of data demand dynamics with switching cost. $\epsilon = 0.5$, $p_2 = 0.35$, $d_2^* = 0.5$, $\gamma_2 = 0.2$. $C_2 = 1/8$ in upper plot and $C_2 = 2/3$ in lower plot.

## 3.6 Numerical Results

In the numerical results, we assume that the congestion costs are given by $g_1(D_1) = D_1/C_1$ and $g_2(D_2) = D_2/C_2$, which capture the congestion externality effects in time-sharing communications networks [39][41]. For the ease of presentation, we consider uniformly distributed $(\theta, d)$, i.e., $f(\theta, d) = 1$ in $\mathcal{U} = \{(\theta, d) \,|\, 0 \le \theta \le 1, 0 \le d \le 1\}$. Note that our analysis also applies to other settings, provided that Assumptions 1–3 specified in Chapter 3.3 are satisfied.

### 3.6.1 Single data plan

First, we illustrate in Fig. 3.1 the oscillation and convergence of the data demand dynamics. The lower plot in Fig. 3.1 shows that the equilibrium data demand $D_2^*$ decreases when the fixed subscription fee $p_2$ increases. Fig. 3.1 verifies that, even for the

Figure 3.4: Single data plan: comparison between the optimal revenue and that yielded by Algorithm I. $\epsilon = 0.5$.

same data plan, different congestion cost functions may result in different convergence behaviors of the data demand dynamics. We plot the continuous-time data demand dynamics specified by (3.20) in Fig. 3.2. It can be seen that the continuous-time and discrete-time data demand dynamics converges to the same equilibrium point. The impacts of switching costs on the data demand dynamics are shown in Fig. 3.3, in which the upper plot indicates that switching costs may make the data demand dynamics converge even though the network capacity is not large enough. We explain this point by noting that, with switching costs, fewer users will not change their subscription decisions and hence the data demand dynamics converges under milder conditions. It can also be seen from the lower plot in Fig. 3.3 that there may exist multiple equilibrium data demand points and the equilibrium, to which the data demand converges, depends on the initial point. Next, we show in Fig. 3.4 that the proposed heuristic Algorithm I can yield a revenue close to the optimum, especially when the network capacity is large. Thus, Algorithm I may be used to find a suboptimal data plan if finding the

(a)                    (b)                   (c)

Figure 3.5: Single data plan: optimal data plan $\mathcal{P}_2 = (p_2, d_2^*, \gamma_2)$ versus capacity $C_2$. $\epsilon = 0.5$.

optimal one is prohibitive. We also plot the optimal data plans[12] $\mathcal{P}_2 = \{p_2, d_2^*, \gamma_2\}$ in Fig. 3.5 under different network capacities. In Fig. 3.6, we show the WSP's profit versus its deployed network capacity under different capacity costs. It indicates that if the average capacity cost is smaller (e.g., the network's lifespan is long and/or the deployment cost is small), then the WSP needs to enlarge its investment in the network capacity. With a larger network capacity, the congestion effects will be reduced and the WSP can attract more users (hence, more revenue) to subscribe to its service.

### 3.6.2 Two data plans

Convergence and oscillation of the data demand dynamics in a wireless market with two data plans are illustrated in Fig. 3.7. As intuitively expected and reflected in Proposition 7, a more stringent requirement on the congestion costs (i.e., the network capacities) is imposed to guarantee the convergence of the data demand dynamics with two data plans, compared to a market with only one data plan. Thus, even though a certain network capacity may guarantee the convergence of data demand dynamics with one data plan, it does not necessarily guarantee the convergence with two data plans. Next, we show in Fig. 3.8 the profits under various network capacities. To maximize the profit, the WSP needs to increase the network capacity for its capped

---

[12]The optimal data plans are obtained by exhaustive search over all the possible data plans.

Figure 3.6: Single data plan: optimal profit versus capacity. $\epsilon = 0.5$.



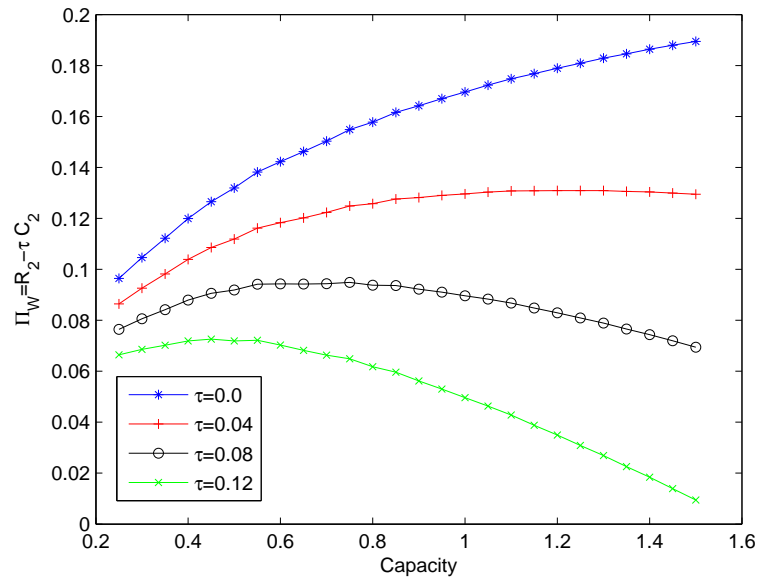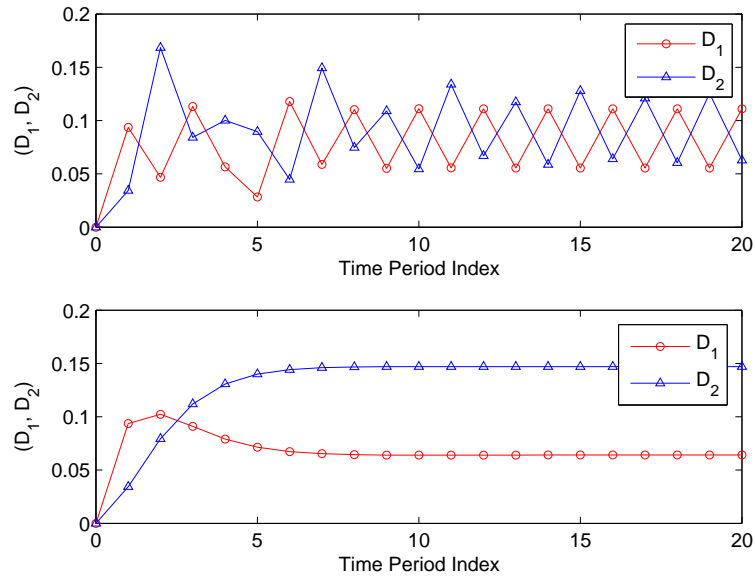Figure 3.7: Two data plans: oscillation and convergence of data demand dynamics. $\epsilon = 0.5$, $p_1 = 0.5$, $p_2 = 0.35$, $\gamma_2 = 0.3$. $C_1 = C_2 = 1/3$ in upper plot and $C_1 = C_2 = 2.0$ in lower plot.

data plan while reducing the network capacity for its unlimited data plan. This can be explained as follows: when an unlimited data plan is offered, subscribers with high data demand will cause excessive congestion costs for the other subscribers, reducing the profitability of the unlimited data plan. This also coincides with the current trend that some WSPs have discontinued unlimited data plans [52]. Finally, we show the competition between two WSPs in Fig. 3.9. It can be seen that if the two WSPs choose their optimal data plans independently in response to the competitor's data plan, then the competition will quickly lead to an equilibrium. This means that, given the long-term capacity investment, the WSPs only adjust their data plans a few times before reaching an equilibrium.[13] Fig. 3.9 also shows that if the capacity investment by the WSP offering the unlimited data plan decreases, the corresponding revenue will be reduced, whereas its competitor's revenue will significantly increase. This is because with a decreased capacity for the unlimited data plan, the resulting congestion cost will increase significantly (due to heavy users) and thus, many users will switch to the usage-based data plan. Note that the two WSPs also need to compete against each other by strategically choosing their long-term network capacities. The result is similar to Fig. 3.9 and hence, is omitted here for brevity.

## 3.7 Conclusion

In this chapter, we considered a wireless communications market where one monopolistic WSP serves a large number of users. The users' data demand dynamics, the WSP's data plan decision and network capacity decision were studied. In our analysis, the users' heterogeneity in terms of their benefits and data demand, as well as the network congestion costs, were explicitly taken into consideration. For the user's data demand dynamics, we showed that: (1) the existence of an equilibrium data demand

---

[13]Similar results are also observed for other simulation settings, although in some (rare) cases the data plan competition between the two WSPs does not converge.

Figure 3.8: Two data plans: optimal profit. $\epsilon = 0.5, \tau = 0.02$.



Figure 3.9: Two data plans: iteration of revenues in a duopoly market. $\epsilon = 0.5$.

is independent of the congestion cost, although for certain data plans, there may not exist any equilibrium data demand if two data plans are both offered in the market; (2)

in order to guarantee the convergence of data demand dynamics, the congestion costs should not increase too rapidly when the aggregate data demand increases, implying that the WSP needs to deploy a large network capacity to support the users' demand. We also proposed a heuristic algorithm that can achieve a close-to-optimal equilibrium revenue if only one data plan is offered by the WSP. For general cases, the WSP's data plan decision and network capacity decision were formalized and solved numerically to maximize the WSP's profit. Finally, we conducted extensive simulations to verify our analysis. Numerical results indicate that to maximize the profit, the WSP should increase the network capacity for its capped data plan while decreasing the network capacity for its unlimited data plan.

# CHAPTER 4

# Maximizing Profit on User-Generated Content Platforms

User-generated content, such as blogs and self-made videos, has becoming a key component in emerging social media. In this chapter, we consider a user-generated content platform monetized through advertising and managed by an intermediary. To maximize the intermediary's profit given rational participants (i.e., content producers and content viewers), we propose a payment scheme in which the intermediary either taxes or subsidizes a content producer an amount of money proportional to the number of views of the producer's content. First, we use a model with a representative content viewer to determine how the content viewers' attention is allocated across available content by solving a utility maximization problem. Then, by modeling the content producers as self-interested agents making independent production decisions, we show that there exists a unique equilibrium in the content production stage, and propose a best-response dynamics to model the decision-making process and to reach the equilibrium. Next, we study the optimal payment scheme (i.e., the payment scheme maximizing the intermediary's profit) that the intermediary chooses taking into account the decisions made by the representative content viewer and the content producers. In particular, by considering the well-known quality-adjusted Dixit-Stiglitz utility function for the representative content viewer, we derive explicitly the optimal payment per content view and characterize analytical conditions under which the intermediary should tax or subsidize the content producers. Finally, we generalize the analysis by

87

Table 4.1: Comparison Between Different Mechanisms

| Mechanism/Scheme | Purpose | Tax | Subsidize | Application Examples |
|---|---|---|---|---|
| Removing low-quality content [58] | Incentivize high-quality content | N/A | N/A | User-generated content |
| Review-based scoring rule [59] | Encourage early responders and high-quality answers | N/A | N/A | Q&A forums |
| Virtual reward [60] | Maximize task competition probability | N/A | Users (virtual currency) | Social networks |
| Payment transfer [61] | Maximize social welfare | Downloaders | Uploaders | P2P networks |
| Pricing transmitters [62] | Maximize network utility | Transmitters | N/A | Wireless relay networks |
| Pricing content providers [99] | Maximize profit | Content providers | Price-sensitive end users | Communications markets |
| Pricing consumers [65] | Maximize profit | Consumers | N/A | Online games & dating |
| **Proposed payment scheme** | Maximize profit | Content producers | Content producers | User-generated content |

considering heterogeneity in terms of production costs among the content producers.

## 4.1 Introduction

As the Internet has been penetrating every aspect of our lives, we have witnessed a significant expansion of online user-generated content platforms during the past decade. Well-known examples include YouTube, Facebook, Twitter, and Yahoo! Answers (an online community where people share knowledge). A key characteristic of these platforms is that the content can be viewed for free by the users and the content producers are not obliged to produce content on the platform. In addition, advertising accounts for a major source of platforms' revenue. To enhance a platform's profitability, it is vital for the platform owner, which we refer to as an *intermediary* in this chapter, to either attract more content views or develop more efficient advertising algorithms showing more relevant advertisement to the content viewers.

In this chapter, we consider a user-generated content platform and propose a payment scheme in which the intermediary can either tax or subsidize the content producers to maximize its profit. On the one hand, the intermediary increases the revenue per content view by taxing, while decreasing the total content available on its platform. On the other hand, the intermediary stimulates more content production by sharing (part of) its advertising revenue with, i.e., subsidizing, the content producers. While

the latter approach has been exercised, either partially or fully, by several online content platforms (e.g., YouTube Partner, Squidoo), we shall show that taxing the content producers, which is relatively less common on the Internet,[1] may also maximize the intermediary's profit.

The scenario we focus on is summarized as follows:

1. The intermediary monetizes the user-generated content platform through advertising, and users can access to the content for free (e.g., YouTube, Yahoo! Answers).[2]

2. Content production is costly. For instance, video content producers may incur costs when shooting video clips and making video files.[3] Content producers are motivated to produce content by both implicit incentives (e.g., social satisfaction) and explicit incentives (e.g., monetary compensation) [70]. In the proposed payment scheme, subsidizing (taxing) provides content producers with an explicit incentive (disincentive) [70].

3. Depending on the payment scheme (i.e., taxing or subsidizing), the intermediary pays a content producer a positive or negative amount proportional to the number of times that its content is viewed. The considered "pay per content view" is a common practice in the Internet industry.

To approach the problem of optimizing the payment scheme, we adopt a leader-follower model (i.e., the intermediary is the leader, followed by the content producers and then by the content viewers) and use backward induction. First, we use the model of a *representative* content viewer, which is a collection of all the individual content viewers, to determine how the content viewers' attention is allocated across a variety of content by solving a utility maximization problem. Then, we study the content

---

[1]The intermediary may tax the content producers for utilizing its resources. As a relevant example, Google Picasa charges its users for storage exceeding the free quota.

[2]In principle, the intermediary may charge the content viewers for viewing the content, which we shall study in our future work.

[3]The production cost also includes the cost (e.g., time cost) incurred in publishing (or uploading) the content on the intermediary's platform.

production decisions made by self-interested content producers. It is shown that there always exists a unique equilibrium point at which no content producer can gain by changing its production decision. We also propose a best-response dynamics to model the content producers' decision process and derive a sufficient condition for its convergence. Next, we formulate the intermediary's profit maximization problem and, by considering the quality-adjusted Dixit-Stiglitz utility function for the representative content viewer, derive a closed-form optimal solution explicitly. We also characterize analytically conditions under which the intermediary should tax or subsidize the content producers. Finally, we generalize our analysis by considering heterogeneity in terms of production costs among the content producers.

The rest of this chapter is organized as follows. Related work is reviewed in Chapter 4.2. Chapter 4.3 describes the model. In Chapter 4.4, we study the decisions made by the content viewers and content producers, and derive the optimal payment maximizing the intermediary's profit. Heterogeneous production costs are studied in Chapter 4.5. Finally, concluding remarks are offered in Chapter 4.6.

## 4.2   Related Works

We summarize in Table 4.1 several mechanisms closely related to ours.

If the intermediary chooses to subsidize the content producers, the proposed payment scheme is essentially an incentive mechanism. Various incentive mechanisms have been proposed recently. For instance, the authors in [58] proposed eliminating or hiding low-quality content to provide content producers with incentives to generate high-quality content. In [59], two scoring rules, the approval-voting scoring rule and the proportional-share scoring rule, were proposed to enable the high-quality answers for online question and answer forums (e.g., Yahoo! Answers). The authors in [60] proposed a (virtual) reward-based incentive mechanism to improve the overall task

completion probability in collaborative social media networks. Pricing-based incentives were proposed to improve social welfare for peer-to-peer networks in [61] and to maximize any system utility in multi-user relay networks in [62], respectively.

If the intermediary taxes the content producers, then the proposed scheme can be classified as market pricing. By considering a general two-sided market, the authors in [63] studied the tradeoffs between the merchant mode and the platform mode, and showed the conditions under which the merchant or platform mode is preferred. Focusing on the Internet markets, [64] revealed that a neutral network is inferior to a non-neutral one in terms of social welfare when the ratio between advertising rates and end user price sensitivity is either too high or too low. The authors in [99] studied the broadband communications market based on a two-sided model, and proposed pricing the content providers to maximize the service provider's profit. In the presence of externalities among the users or consumers, [65] proposed three different pricing schemes, namely, uniform price, differentiated prices and only two different prices, to maximize the monopolist's profit.

## 4.3   Model

Consider an online user-generated content platform managed by a profit-maximizing intermediary. Content on the platform is produced by individual users, which we refer to as *content producers*, and viewed by content viewers. Next, we specify the modeling details of the intermediary, content producers and viewers.

### 4.3.1   Intermediary

It is known that advertising is one of the most prevailing revenue sources in the Internet industry, especially for content platforms such as YouTube and Yahoo Answers. Hence, we consider that the intermediary monetizes its content platform by displaying

contextual advertisement to content viewers. Although multiple charging models (e.g., pay per click, pay per sale, and pay per impression) have been proposed for online advertising, the advertising revenue is in general proportional to the total number of times that the content with advertisement is viewed (i.e., content views) [64]. To increase the advertising revenue, one natural approach is to improve the algorithm based on which the advertisement is displayed to the content viewers, e.g., targeted advertising that enhances the correlation of advertisement to the content such that the advertisement is more likely to be clicked on. On the other hand, the intermediary can derive more advertising revenue by increasing the content views. To do so, we propose that the intermediary provides the content producers with economic incentives (i.e., subsidizing) to produce more content,[4] which in turn attracts more content views. Essentially, this scheme allows the intermediary to share with the content producers (part of) its advertising revenue as an economic incentive, and it has been exercised by several content platforms (e.g., YouTube Partner, Squidoo). Note that, since we focus on the side of content producers, we do not consider subsidizing the content viewers (e.g., providing rewards) to attract more content views. Moreover, the number of content viewers is typically much greater than that of content producers and hence, in practice, it is relatively easier to implement the proposed payment scheme on the side of content producers.

To formally state our model, we denote $\bar{x}$ as the total content views of all the content on the intermediary's platform, and $b \geq 0$ as the (average) advertising profit (i.e., revenue minus cost) that the intermediary can derive per content view. For the convenience of analysis, we assume that $b$ is constant regardless of $\bar{x}$, i.e., the average advertising profit is independent of the content views. The intermediary pays $\theta$ per content view to the respective content producers. For the completeness of analysis, we allow $\theta$ to take negative values, in which case the intermediary taxes the content

---

[4]In this chapter, the intermediary does not differentiate among different qualities, although a producer producing content of higher quality will be compensated more because of higher demand.

producers with $-\theta$ per content view. Practically speaking, negative $\theta$ may correspond to that the intermediary taxes the content producers for utilizing its resources (e.g., bandwidth) or for commission fees if the content producers produce advertisement-type content. In the following analysis, we use the term *payment* (per content view) to refer to $\theta$ wherever applicable, regardless of its positive or negative sign. Neglecting the intermediary's recurring fixed operational cost, we can express the intermediary's profit as

$$\Pi_{\mathcal{I}} = (b - \theta) \cdot \bar{x}. \tag{4.1}$$

While $b$ can be increased by using sophisticated advertising algorithms showing more relevant advertisement, we assume throughout the chapter that $b$ is exogenously determined and fixed, and shall focus on deriving the optimal $\theta$ that maximizes the intermediary's profit. Note that in our current study, we restrict our analysis to uniform $\theta$, and considering more general payment schemes (e.g., different payments per content view for different content producers) is left for future work.

### 4.3.2 Content Producers

As evidenced by the exploding number of YouTube users, a popular user-generated content platform can attract a huge number of content producers. To capture this fact, we use a continuum model and assume that the mass of content producers is normalized to one. Each content producer can produce content of a certain quality while incurring a production cost. We further assume that content producers produce differentiated content, or in other words, no two content producers can produce identical content. Note that the content quality can be different across content producers, although we assume that the production cost is the same for all content producers (we shall relax this assumption in Chapter 4.5). The content quality is represented by a scalar and treated as an internal feature of content (e.g., how fun/informative the content is). Mathematically, we denote $q_i \geq 0$ and $c > 0$ as the quality of content

produced by content producer $i$ and the production cost, respectively. Without causing ambiguity, we occasionally use content $q_i$ to refer to the content with a quality $q_i$. To characterize heterogeneity in the content quality, we assume that the content quality follows a cumulative distribution function (CDF), denoted by $F(q)$, across the unit mass of content producers. In other words, $F(q)$ denotes the number or fraction of content producers whose content has a quality less than or equal to $q \geq 0$.

Millions of users engaging daily in Internet activities such as blogs, for which they receive no monetary rewards, highlight that such content producers may simply derive satisfaction (and hence utility) by attracting the content viewers' attention [58][60][70]. We use the content views to quantify the amount of received attention and assume that the benefit resulting from the content viewers' attention for a content producer is a linear function of its content views. We assume further that each content producer $i$ is self-interested and can strategically make a binary decision: produce or not produce. Denote by $x(q_i) \geq 0$ the number of views for content $q_i$. If content producer $i$ produces content on the intermediary's platform, it can derive a utility expressed as

$$\pi_i = (\theta + s) \cdot x(q_i) - c, \tag{4.2}$$

where $s \geq 0$ is the (social) benefit per content view derived from the content viewers' attention,[5] $\theta$ is the payment per content view determined by the intermediary, and $c$ is the production cost. Content producer $i$ obtains zero utility if it chooses not to produce content. By the assumption of rationality, content producer $i$ chooses to produce content if and only if its utility is non-negative.

In what follows, we assume that the content quality $q$ follows a distribution in a normalized interval $[0, 1]$ and the probability density distribution (PDF) is given by a continuous and positive function $f(q)$ for $q \in [0, 1]$. Scaling the interval $[0, 1]$ to $[0, \bar{q}]$

---

[5]Essentially, $s \geq 0$ converts the content viewer's attention to the content producers' (economic) benefit/utility.

does not affect the analysis, but will only complicate the notations. It is intuitively expected that a content with a higher quality will attract more content views (and yield a higher utility for its content producer, too) than the one with a lower quality. Thus, the production decision of the content producers has a threshold structure. In particular, there exist *marginal* content producers whose content has a quality denoted by $q_m \in [0, 1]$, and those content producers whose content quality is greater (less) than $q_m$ will (not) choose to produce content on the intermediary's platform [67]. We refer to $q_m$ as the marginal content quality. Next, it is worthwhile to provide the following remarks concerning the model of content producers.

*Remark 1:* In our model, a content producer who produces $m \geq 1$ pieces of content is viewed as $m$ content producers, each of whom produces a single content, and the total production cost is $m \cdot c$ (i.e., constant returns to scale [63]).

*Remark 2:* As in [63], we assume that the content producers will incur a predetermined production cost if they choose to produce content. That is, the content producers cannot choose their production costs when producing content. For the ease of presentation and to gain insights as to whether the intermediary should tax or subsidize the content producers, we first consider a homogeneous production cost among the content producers. In Chapter 4.5, we shall generalize the model to consider heterogeneity in the content producers' production costs.

*Remark 3:* If $\theta < -s$, it is clear from (4.2) that no content producers can possibly receive a non-negative utility by producing content on the platform. As a consequence, $\bar{x} = 0$ and the intermediary's profit is zero. On the other hand, if $\theta > b$, then we see from (4.1) that the intermediary can never gain a positive profit. Hence, we exclude these two trivial cases in the remainder of this chapter and focus on the case of $-s \leq \theta \leq b$ unless otherwise stated.

### 4.3.3 Content Viewers

Despite that the content viewers are diverse in terms of preferences towards the content, the aggregate content viewing decisions of all the content viewers can be conveniently characterized by the decision of a representative content viewer. Thus, we adopt the widely-used representative agent model to determine how the total content views are allocated across a variety of content [107]. Specifically, the representative content viewer optimally allocates its total content views, denoted by $T$, across the available content to maximize its utility. Note that $T$ can be interpreted as the size of the representative content viewer or the market size. On the Internet, it is quite common that multiple content platforms offer similar services and the content viewers have access to the content on any of these platforms. Focusing on the intermediary's optimal payment decision, we do not consider the details of how the content is produced on the other platforms. Instead, we can assume that the mass of content available on the other platforms is $n_a \geq 0$ and the content quality follows a certain CDF $\tilde{F}(q)$ with support $q \in [q_l, q_h]$, where $0 \leq q_l < q_h$ are the lowest and highest content quality on the other platforms, respectively. For the convenience of notation, throughout the chapter, we alternatively represent the content on the other platforms using a unit mass of content with an *aggregate* quality of $q_a$, without affecting the analysis. Note that $q_a$ is a function of $n_a \geq 0$, $\tilde{F}(q)$ and the utility function of the representative content viewer. In particular, given a uniform distribution of content quality on the other platforms and the quality-adjusted Dixit-Stiglitz utility for the representative content viewer (which we shall define later), we can readily obtain

$$
q_a = \left[ \frac{n_a \left( q_h^{\sigma+1} - q_l^{\sigma+1} \right)}{1 + \sigma} \right]^{\frac{1}{\sigma}},
\tag{4.3}
$$

where $\sigma > 1$ measures the content substitutability. Recalling that $q_m \in [0, 1]$ is the marginal content quality above which the content producers choose to produce con-

tent on the intermediary's platform, we write the representative content viewer's utility function as $U(x(q), x_a \mid q_m, q_a)$, where $x(q)$ denotes the content view for content $q \in [q_m, 1]$ and $x_a$ is the content view allocated to the aggregate content $q_a$ on the other content platforms. In our model, the intermediary's content platform is a monopolist in the market and the content producers under consideration, if they choose to produce content, can only produce content on the intermediary's platform. Thus, $x_a$ is essentially interpreted as "outside activity" of the content viewers. Note that $x(q)$ can be rewritten as $x(q \mid q_m, q_a)$, although we use the succinct notation $x(q)$ throughout the chapter whenever applicable. If $q_m$ increases (decreases), there will be less (more) content on the intermediary's platform. Because of the continuum model, we allow $x(q)$ and $x_a$ to take non-integer values, and $x(q)$ actually represents the content view *density* allocated to a continuum of content with quality $q \in [q_m, 1]$, i.e., $x(q)$ is the content view that an *individual* content producer with a content quality of $q$ receives. Next, we formulate the utility maximization problem for the representative content viewer as follows

$$
\begin{aligned}
\max_{x(q) \geq 0, x_a \geq 0} \ & U(x(q), x_a \mid q_m, q_a), \\
s.t., \quad & \int_{q_m}^{1} x(q) dF(q) + x_a \leq T,
\end{aligned}
\tag{4.4}
$$

where $F(q)$ is the CDF of content quality on the intermediary's content platform. It is worth noting that an implicity assumption underlying the problem (4.4) is that the aggregate quality of the content on the other platforms is independent of the intermediary's payment decision and other variables in the model such as $q_m$, $x(q)$, $x_a$. This can be justified by noting that there are many content platforms on the Internet and changes on one content platform have a negligible impact on the other platforms. Before performing further analysis, we assume that the following properties are satisfied by the utility function $U(x(q), x_a \mid q_m, q_a)$.

*Property 1 (Diminishing marginal utility):* $U(x(q), x_a \mid q_m, q_a)$ is increasing and

strictly (jointly) concave in $x(q)$ and $x_a$, for $q \in [0,1]$.

*Property 2 (Preference towards diversified content):* $\max_{x(q) \geq 0, x_a \geq 0} U(x(q), x_a \,|\, q_m, q_a)$ is decreasing in $q_m \in [0,1]$.

*Property 3 (Negative externalities):* Denote by $x^*(q \,|\, q_m, q_a)$, for $q \in [0,1]$, the optimal solution to (4.4).[6] If content $q$ is produced, then $x^*(q \,|\, q_m, q_a)$ is positive. Moreover, it is continuous and strictly increasing in $q_m \in [0,1]$, increasing in $q \in [0,1]$, and decreasing in $q_a$ for $q_a \in [0, \infty)$. In particular, $x^*(0 \,|\, q_m, q_a) = 0$ for all $q_m \in [0,1]$ and $q_a \geq 0$.

*Property 4 (More content leading to more content views):* $\bar{x} = \int_{q_m}^1 x^*(q \,|\, q_m, q_a) dF(q)$ is decreasing in $q_m \in [0,1]$.

We briefly discuss the above properties. Property 1 captures the effects of diminishing marginal utility when the representative content viewer views more content. Property 2 models the phenomenon that content viewers will typically benefit from the participation of content producers on the platform. This is particularly true for online content platforms, where the content viewers prefer to view a diversified bundle of content. Thus, when $q_m \in [0,1]$ increases, i.e., fewer content producers produce content, the representative content viewer's (maximum) utility decreases. Property 3 reflects the "crowding effects", i.e., lower $q_m$ or more content production increases competition among the content producers. Specifically, an individual content producer will attract a less content view if more content producers choose to produce content on the platform or the aggregate content quality on the other platforms is higher. The last property ensures that more content views are devoted to the intermediary's platform if there is more content available on the platform.

As a concrete example satisfying Properties 1–4, we use a *quality-adjusted* version

---

[6]The existence of $x^*(q \,|\, q_m, q_a)$ can be established based on the strict concavity of the utility function. Due to the continuum model, $\tilde{x}(q \,|\, q_m, q_a) = x^*(q \,|\, q_m, q_a)$ *almost everywhere* for $q \in [0,1]$ is also optimal in maximizing the utility function. For simplicity, we treat such $\tilde{x}(q \,|\, q_m, q_a)$ the same as $x^*(q \,|\, q_m, q_a)$. The treatment does not change our analysis, except that it affects the decisions of a negligible mass of content producers.

of the well-known Dixit-Stiglitz utility function [68][107], defined as below[7]

$$U(x(q), x_a \mid q_m, q_a) = \left[ \int_{q_m}^1 q \cdot x(q)^{\frac{\sigma-1}{\sigma}} dF(q) + q_a \cdot x_a^{\frac{\sigma-1}{\sigma}} \right]^{\frac{\sigma}{\sigma-1}}, \qquad (4.5)$$

where $\sigma > 1$ measures the elasticity of substitution between different content. In the extreme case, the content is perfectly substitutable when $\sigma = \infty$ [68].

## 4.4 Profit Maximization on Content Platforms

In this part, based on the model described above and the proposed payment scheme, we study the problem of optimizing payment in the presence of self-interested content producers and content viewers. The timing of decision making can be described as follows.

**Stage 1 (Payment Scheme Decision):** The intermediary announces the value of $\theta$ to the content producers.

**Stage 2 (Production Decision):** Given $\theta$, each content producer makes a binary decision as to whether or not to produce content on the platform.

**Stage 3 (Viewing Decision):** Given the available content, the content viewers, consolidated as a representative content viewer, optimally allocate the content views to maximize utility subject to a total content view constraint.

From the described timing, we see that the intermediary can be regarded as the leader, followed by the content producers and then by the content viewers. Thus, in order to identify the optimal $\theta$, the intermediary needs to first know how the content producers and viewers respond to its choice of $\theta$. Therefore, we proceed with our

---

[7]If we model the quality of $n_a$ pieces of content on the other platforms by using a certain CDF $\tilde{F}(q)$ for $q \in [q_l, q_h]$, where $0 \leq q_l < q_h$ are the lowest and highest content quality on the other platforms, respectively, then the Dixit-Stiglitz utility function in (4.5) becomes $U(x(q) \mid q_m, q_l, q_h) = \left[ \int_{q_m}^1 q \cdot x(q)^{\frac{\sigma-1}{\sigma}} dF(q) + n_a \int_{q_l}^{q_h} q \cdot x(q)^{\frac{\sigma-1}{\sigma}} d\tilde{F}(q) \right]^{\frac{\sigma}{\sigma-1}}$ and the analysis remains the same with an appropriate choice of $q_a$.

analysis using backward induction.

### 4.4.1 Optimal Content Viewing

It follows from the strict concavity of $U(x(q), x_a \mid q_m, q_a)$, specified in Property 1, that there exists a unique optimal solution, denoted by $x^*(q)$ and $x_a^*$ for $q \in [0, 1]$, to the utility maximization problem in (4.4), although it is not possible to obtain a closed-form expression without specifying $U(x(q), x_a \mid q_m, q_a)$. By considering the quality-adjusted Dixit-Stiglitz utility defined in (4.5) and uniform distribution of the content quality,[8] we can obtain explicitly the closed-form solution as follows

$$x^*(q) = \frac{T(\sigma+1)q^\sigma}{(\sigma+1) \cdot q_a^\sigma + (1-q_m^{\sigma+1})}, \tag{4.6}$$

for $q \in [q_m, 1]$, $x^*(q) = 0$ for $q \in [0, q_m)$, and $x_a^* = \frac{T(\sigma+1)q_a^\sigma}{(\sigma+1) \cdot q_a^\sigma + (1-q_m^{\sigma+1})}$. The details of deriving (4.6) are omitted for brevity. After plugging $x^*(q)$ and $x_a^*$ into (4.5), the maximum utility derived by the representative content viewer is given by

$$U^*(x^*(q), x_a^*) = T \left[ q_a^\sigma + \frac{1 - q_m^{\sigma+1}}{\sigma+1} \right]^{\frac{1}{\sigma-1}}, \tag{4.7}$$

which is decreasing in $q_m \in [0, 1]$.

### 4.4.2 Equilibrium Content Production

Due to rationality, content producers will not choose to produce content if they cannot obtain non-negative utilities. Essentially, interaction among the content producers can be formalized as a non-cooperative game with an infinite number of players, the solution to which is (Nash) equilibrium. At an equilibrium, if any, no content producers

---

[8]The uniform distribution has been widely applied to model the diversity of various factors, such as opportunity cost [63] and valuation of quality-of-service [66].

can gain more benefits by deviating from their decisions. In other words, the fraction of content producers choosing to produce content on the intermediary's platform does not change at the equilibrium, or equivalently, the marginal content quality $q_m \in [0, 1]$ becomes invariant. Next, we study the equilibrium content production by specifying the equilibrium marginal content quality denoted by $q_m^*$.

If $q_m^* = 1$, then no (or a zero mass of) content producers can receive a non-negative utility by producing content on the platform. This implies that, with $q_m^* = 1$, we have $x^*(1|1, q_a) \cdot (\theta + s) - c \leq 0$. If there are some content producers choosing to produce content at the equilibrium (i.e., $q_m^* \in [0, 1)$), then according to the definition of marginal content producers, we have $x^*(q_m^*|q_m^*, q_a) \cdot (\theta + s) - c = 0$. Hence, we can show that $q_m^* \in [0, 1]$ satisfies

$$q_m^* \triangleq Q(q_m^*) = \arg\min_q \left\{ q \in [0, 1] : \ x^*(q \,|\, q_m^*, q_a) \cdot (\theta + s) - c \geq 0 \right\}, \quad (4.8)$$

where $[x]_0^1 = \min\{1, \max\{0, x\}\}$. Thus, an equilibrium point of content production exists if and only if the mapping $Q(q_m^*)$, defined in (4.8), has a fixed point. Next, we formally define the equilibrium marginal content quality in terms of $q_m^*$ as below.

*Definition 1:* $q_m^*$ is an *equilibrium* marginal content quality if it satisfies $q_m^* = Q(q_m^*)$.

We establish the existence and uniqueness of an equilibrium marginal content quality in Theorem 3.

**Theorem 3.** *For any $\theta \in [-s, b]$, there exists a unique equilibrium $q_m^* \in (0, 1]$ in the production decision stage. Moreover, $q_m^*$ satisfies*

$$\begin{cases} q_m^* = 1, & \text{if } x^*(1 \,|\, 1, q_a) \cdot (\theta + s) \leq c, \\ q_m^* \in (0, 1), & \text{otherwise,} \end{cases} \quad (4.9)$$

101

*where $x^*(1\,|\,1, q_a)$ is obtained by solving (4.4) with $q_m \to 1$.*[9]    $\square$

*Proof.* We prove the theorem by considering the following two cases.

Case 1: If $x^*(1\,|\,1, q_a) \cdot (\theta + s) \leq c$, then it can be shown that for any value of $q_m \in [0, 1]$, the utility obtained by content producer $i$ whose content quality is 1 satisfies

$$\pi_i = x^*(1\,|\,q_m, q_a) \cdot (\theta + s) - c \leq x^*(1\,|\,1, q_a) \cdot (\theta + s) - c \leq 0, \qquad (4.10)$$

where the first inequality follows from Property 3 specified in Chapter 4.3 that $x^*(1\,|\,q_m, q_a) \leq x^*(1\,|\,1, q_a)$. Thus, only $q_m^* = 1$ satisfies Definition 1 and the equilibrium is unique in this case.

Case 2: If $x^*(1\,|\,1, q_a) \cdot (\theta + s) > c$, we first show by contradiction that $q_m^*$, if any, must be strictly less than 1. Suppose that $q_m^* = 1$ and no content is produced at the equilibrium. Then, content producer $i$ whose content quality is 1 can receive a strictly positive utility

$$\pi_i = x^*(1\,|\,q_m^*, q_a) \cdot (\theta + s) - c = x^*(1\,|\,1, q_a) \cdot (\theta + s) - c > 0, \qquad (4.11)$$

which means that content producer $i$ can receive a positive utility by producing content and contradicts our assumption. Thus, $q_m^*$ must be strictly less than 1. On the other hand, by Property 3, we have $x^*(0\,|\,q_m, q_a) = 0$, for any $q_m \in [0, 1]$, and hence content producers whose content quality is zero will not choose to produce content. Therefore, at the equilibrium, $q_m^*$ must be strictly positive.

Next, we prove the uniqueness of $q_m^*$ by constructing an auxiliary function $\bar{Q}(q_m) = Q(q_m) - q_m$, where $Q(q_m)$ is defined in (4.8) is continuous in $q_m \in [0, 1]$ as per Property

---

[9]When $q_m \to 1$, only a negligible fraction of content producers choose to produce content on the intermediary's content platform.

3. Thus, $q_m^*$ is an equilibrium point of the content production if and only if $\bar{Q}(q_m^*) = 0$ (i.e., a fixed point of $Q(q_m)$). Based on Property 3, it can be seen that, when $q_m$ increases, $Q(q_m)$ decreases. Thus, the auxiliary function $\bar{Q}(q_m) = Q(q_m) - q_m$ is continuous and strictly decreasing in $q_m \in [0, 1]$. It can be further shown that $\bar{Q}(1) < 0$ while $\bar{Q}(0) > 0$. Thus, by applying the intermediate value theorem, there must exist a unique $q_m^* \in (0, 1)$ such that $\bar{Q}(q_m^*) = 0$. Therefore, the uniqueness of the equilibrium point in the content production stage is proved. ∎

Theorem 3 guarantees the existence of a unique equilibrium point and shows that if the the content producer with the highest quality cannot obtain a positive utility (due to high production cost, taxing or low subsidizing from the intermediary), then no content producers choose to produce content on the intermediary's content platform at the equilibrium. For notational convenience, we denote the value of $\theta$ that satisfies $x^*(1 \,|\, 1, q_a) \cdot (\theta + s) = c$ by

$$\underline{\theta} = \frac{c}{x^*(1 \,|\, 1, q_a)} - s. \tag{4.12}$$

Then, it follows from Theorem 3 that the intermediary can gain a positive profit if and only if $\theta \in (\underline{\theta}, b]$. Nevertheless, if $\underline{\theta} \geq b$, then the intermediary's profit is always zero. Hence, we assume $\underline{\theta} < b$ throughout the chapter. Based on the uniqueness of $q_m^*$ for any $\theta \in [-s, b]$, we can express $q_m^* = q_m^*(\theta)$ as a function of $\theta \in [-s, b]$. While there exists no simple closed-form expression of $q_m^*(\theta)$ in general, it can be easily shown that $q_m^*(\theta) \in (0, 1)$ is strictly decreasing in $\theta \in (\underline{\theta}, b]$ and $q_m^*(\theta) = 1$ for $\theta \in [-s, \underline{\theta}]$. In practice, the content producers do not have complete information regarding each other and hence, they may not make directly the decisions that strikes an equilibrium. In such a scenario, the content producers may use an adjustment process to update their decisions based on limited information.

A natural and well-studied approach to modeling an adjustment process is the best-

response dynamics, in which each decision maker chooses the best action in response to the decisions made by the others. In this chapter, we consider the best-response dynamics based on naive (or static) expectation. Specifically, at the end of time $t = 1, 2, 3 \cdots$, content producer $i$ assumes that the decisions made by the other content producers at time $t + 1$ remain the same as those at time $t$, and expects $x_{t+1}(q_i) = x^*(q \mid q_{m,t}, q_a)$, where $x^*(q \mid q_{m,t}, q_a)$ is the solution to (4.4) and $q_{m,t} \in [0, 1]$ is the marginal content quality at time $t$. Note that a content producer with a content quality less than $q_{m,t}$ may also choose to produce content at time $t + 1$, if it believes that there is not much high quality content on the platform (i.e., $q_{m,t} \in [0, 1]$ is large) and it can receive a non-negative utility. Similar decision processes have been adopted in the existing literature (e.g., [66] and references therein). The best-response decision model implies that the sequence $q_{m,t}$, for $t = 0, 1, 2 \cdots$, evolves as follows[10]

$$q_{m,t+1} = Q(q_{m,t}), \tag{4.13}$$

where $Q(\cdot)$ is defined in (4.8). Essentially, the dynamics in (4.13) is a fixed point iteration for $Q(\cdot)$ and it converges regardless of the initial point if $|Q'(q)| < 1$ for $q \in [0, 1]$ [69]. Following the contraction mapping theorem, we can easily further specify a sufficient convergence condition, and the details are omitted here for brevity. It should be noted that, by considering the dynamics specified by (4.13), we implicitly assume that the content produced in the previous periods has little value and will not significantly affect the content views in the current period (e.g., news content platform). Moreover, the dynamics specified by (4.13) requires that all the content producers update production decisions at the end of each time period. In practice, if only a fraction $\epsilon \in (0, 1]$ of the content producers make decisions each time, then the sequence becomes $q_{m,t+1} = (1 - \epsilon)q_{m,t} + \epsilon Q(q_{m,t})$ without affecting the equilibrium analysis while the convergence is slowed down.

---

[10]If there exists no $q \in [0, 1]$ such that $x(q \mid q_{m,t}, q_a) \cdot (\theta + s) - c \geq 0$, then we set $q_{m,t+1} = 1$.

### 4.4.3 Optimal Price

Based on decisions made by the content viewers and content producers, we study the optimal payment $\theta$ that maximizes the intermediary's steady-state profit (i.e., profit obtained when the content production decision stage reaches the equilibrium). Mathematically, we formalize the profit maximization problem as

$$\theta^* = \arg \max_{\theta \in [\underline{\theta}, b]} (b - \theta) \cdot \bar{x}, \tag{4.14}$$

where $\bar{x} = \int_{q_m^*}^1 x^*(q \,|\, q_m^*, q_a) dF(q)$. The decision interval is shrunk to $[\underline{\theta}, b]$, since $\theta \in [-s, \underline{\theta})$ always results in a zero profit for the intermediary, where $\underline{\theta}$ is defined in (4.12). By Property 4 stated in Chapter 4.3, $\bar{x} = \int_{q_m^*}^1 x^*(q) dF(q)$ is decreasing in $q_m^* \in [0, 1]$. Then, recalling that $q_m^*(\theta)$ is strictly decreasing in $\theta \in [\underline{\theta}, b]$, we can see $\bar{x}$ is increasing in $\theta \in [\underline{\theta}, b]$. In other words, increasing $\theta \in [\underline{\theta}, b]$ will encourage more content producers to produce content on the intermediary's platform and hence attract more content views from the content viewers.

Although we can numerically solve the profit maximization problem (4.14), it is rather challenging, if not impossible, to explicitly determine the optimal payment $\theta^*$ without specifying the utility function $U(x(q), x_a \,|\, q_m, q_a)$ or further restrictions on $\bar{x}$. If the profit function in (4.1) is strictly concave in $\theta$, then there exists a unique optimal payment $\theta^* \in [\underline{\theta}, b]$ maximizing the intermediary's profit and satisfying the first-order optimality condition

$$-\bar{x}(\theta^*) + (b - \theta^*) \frac{\partial \bar{x}}{\partial \theta}\Big|_{\theta = \theta^*} = 0. \tag{4.15}$$

Moreover, if the first-order partial derivative of (4.1) with respect to $\theta$ evaluated at $\theta = 0$ is negative (positive), then the optimal payment $\theta^*$ is negative (positive), i.e., the intermediary should tax (subsidize) the content producers to maximize its profit. This

result can be interpreted as follows. If $-\bar{x}(\theta = 0) + b \cdot \frac{\partial \bar{x}}{\partial \theta}|_{\theta=0} < 0 \implies \bar{x}(\theta = 0) > b \cdot \frac{\partial \bar{x}}{\partial \theta}|_{\theta=0}$, it implies that there is already a sufficient amount of content available on the intermediary's platform even when the intermediary does not subsidize the content producers. Otherwise, the intermediary should subsidize the content producers such that there is more content produced, receiving more content views. Therefore, we see that the amount of content views when $\theta = 0$ is critical in determining whether the intermediary should tax or subsidize the content producers.

In the following analysis, to gain insights and explicitly derive the optimal payment $\theta^*$, we consider the quality-adjusted Dixit-Stiglitz utility and uniform distribution of content quality. A closed-form optimal payment $\theta^* \in [\frac{c \cdot q_a^\sigma}{T} - s, b]$ is explicitly obtained and shown in Theorem 4.

**Theorem 4.** *Suppose that $U(x(q), x_a \mid q_m, q_a)$ is given by the Dixit-Stiglitz utility function in (4.5) and the content quality $q$ is uniformly distributed on $[0, 1]$. The unique optimal payment $\theta^* \in [\frac{c \cdot q_a^\sigma}{T} - s, b]$ that maximizes the intermediary's profit is given by*

$$\theta^* = \frac{c\left[(\sigma + 1) \cdot q_a^\sigma + 1 - z^{\sigma+1}\right]}{T(\sigma + 1) \cdot z^\sigma} - s, \tag{4.16}$$

*where $z \in [q_m^*(b), 1]$ is the unique root of the equation*

$$-\frac{T \cdot q_a^\sigma \cdot (b + s)}{((\sigma + 1) \cdot q_a^\sigma + 1 - z^{\sigma+1})^2} + \frac{c}{(\sigma + 1)^3} \cdot \frac{\sigma + z^{\sigma+1}}{z^{2\sigma+1}} = 0. \tag{4.17}$$

*Proof.* We see from Theorem 3 that $q_m^* = Q(1) = 1$ and $q_m^* \in (0, 1)$, when $\theta = \frac{c \cdot q_a^\sigma}{T} - s$ and $\theta = b$, respectively. Hence, we can express

$$q_m^* = Q(q_m^*) = \left\{ \frac{c\left[(\sigma + 1) \cdot q_a^\sigma + 1 - (q_m^*)^{\sigma+1}\right]}{T(\sigma + 1)(\theta + s)} \right\}^{\frac{1}{\sigma}} \tag{4.18}$$

when $\theta^* \in [\frac{c \cdot q_a^\sigma}{T} - s, b]$. Then, by rewriting (4.18), $\theta$ can be expressed in terms of $q_m^*$ as

$$\theta = \frac{c\left[(\sigma + 1) \cdot q_a^\sigma + 1 - (q_m^*)^{\sigma+1}\right]}{T(\sigma + 1)(q_m^*)^\sigma} - s. \qquad (4.19)$$

Thus, by replacing the optimization variable $\theta$ with $q_m^*$, we can reformulate the intermediary's profit maximization problem in (4.14) as

$$\max_{q_m^* \in [q_m^*(b), 1]} (b + s - \frac{c\left[(\sigma + 1) \cdot q_a^\sigma + 1 - (q_m^*)^{\sigma+1}\right]}{T(\sigma + 1)(q_m^*)^\sigma}) \cdot \bar{x}, \qquad (4.20)$$

where $q_m^*(b)$ is the equilibrium point of the content production when $\theta = b$. By plugging $\bar{x} = \int_{q_m^*}^{1} x^*(q)dq = T \cdot \frac{1 - (q_m^*)^{\sigma+1}}{(\sigma+1) \cdot q_a^\sigma + 1 - (q_m^*)^{\sigma+1}}$, the objective function in (4.20) can be written as

$$(b + s) \cdot T - \frac{(b + s) \cdot T \cdot (\sigma + 1) \cdot q_a^\sigma}{(\sigma + 1) \cdot q_a^\sigma + 1 - (q_m^*)^{\sigma+1}} - \frac{c[1 - (q_m^*)^{\sigma+1}]}{(\sigma + 1)(q_m^*)^\sigma}, \qquad (4.21)$$

which we can prove is a strictly concave function of $q_m^* \in [q_m^*(b), 1]$. Therefore, there exists a unique solution, denoted by $q_m^{**}$, which maximizes (4.21), and correspondingly, the optimal value of $\theta$ that maximizes (4.14) is also unique in $[\frac{c \cdot q_a^\sigma}{T} - s, b]$.

By applying Theorem 3, we have $q_m^* \in (0, 1)$ if $\theta > \frac{c \cdot q_a^\sigma}{T} - s$, i.e., a positive fraction of content producers produce content and the total content views $\bar{x}$ is positive. Thus, if $\frac{c \cdot q_a^\sigma}{T} - s < b$, there exists a $\theta \in [\frac{c \cdot q_a^\sigma}{T} - s, b]$ such that the intermediary's profit $(b - \theta) \cdot \bar{x}$ is strictly positive.

We see that the first-order derivative of (4.21) with respect to $q_m^*$ can be obtained as

$$-\frac{T \cdot q_a^\sigma \cdot (b + s)(\sigma + 1)^2 (q_m^*)^\sigma}{[(\sigma + 1) \cdot q_a^\sigma + 1 - (q_m^*)^{\sigma+1}]^2} + \frac{c}{\sigma + 1} \cdot \left[\frac{\sigma(q_m^*)^{\sigma-1}}{(q_m^*)^{2\sigma}} + 1\right]. \qquad (4.22)$$

Due to the strict concavity of (4.21), its first-order derivative in (4.22) is strictly de-

creasing in $q_m^* \in [q_m^*(b), 1]$. We have already shown that the intermediary's profit is zero if $\theta = b$ or $\theta = \frac{c \cdot q_a^\sigma}{T} - s$ and positive if $\theta \in (\frac{c \cdot q_a^\sigma}{T} - s, b)$. Thus, there must exist a unique value of $\theta^* \in [\frac{c \cdot q_a^\sigma}{T} - s, b]$ such that (4.22) is zero and the intermediary's profit is maximized. By letting $z = q_m^*$ and dividing (4.22) by $(\sigma + 1)^2 \cdot z^\sigma$, we have proved Theorem 4. ∎

Having derived the optimal payment $\theta^*$, we next analyze the sign of the optimal payment in Proposition 8. Such analysis is useful in understanding the impacts of various factors on the intermediary's decision of taxing or subsidizing.

**Proposition 8.** *Suppose that $U(x(q), x_a \mid q_m, q_a)$ is given by the Dixit-Stiglitz utility function in (4.5) and the content quality $q$ is uniformly distributed in $[0, 1]$. The optimal payment $\theta^* \in [\frac{c \cdot q_a^\sigma}{T} - s, b]$ that maximizes the intermediary's profit satisfies*

$$
\begin{cases}
\theta^* \in (0, b), & if \, \Delta < \dfrac{c(\sigma + 1) \cdot q_a^\sigma \cdot (b + s)}{s^2 T}, \\[2mm]
\theta^* = 0, & if \, \Delta = \dfrac{c(\sigma + 1) \cdot q_a^\sigma \cdot (b + s)}{s^2 T}, \\[2mm]
\theta^* \in (\dfrac{c \cdot q_a^\sigma}{T} - s, 0), & if \, \Delta > \dfrac{c(\sigma + 1) \cdot q_a^\sigma \cdot (b + s)}{s^2 T},
\end{cases}
\tag{4.23}
$$

*where $\Delta = \frac{\sigma}{q_m^*(0)} + [q_m^*(0)]^\sigma$, in which $q_m^*(0)$ is the equilibrium point of content production when the intermediary chooses $\theta = 0$.*

*Proof.* Assuming that $\theta^* = 0$, we obtain from (4.16) that

$$
(\sigma + 1) \cdot q_a^\sigma + 1 - [q_m^*(0)]^{\sigma+1} = \frac{sT(\sigma + 1)[q_m^*(0)]^\sigma}{c}.
\tag{4.24}
$$

Then, by plugging (4.24) into (4.22), we can rewrite (4.22) as

$$
\frac{c}{(\sigma + 1)[q_m^*(0)]^\sigma} \left[ -\frac{c(\sigma + 1) \cdot q_a^\sigma \cdot (b + s)}{s^2 T} + \frac{\sigma}{q_m^*(0)} + [q_m^*(0)]^\sigma \right],
\tag{4.25}
$$

(a) Impacts of $q_a$. $T = 10$, $c = 1.0$.

(b) Impacts of $c$. $T = 10$, $q_a = 1.5$.

(c) Impacts of $T$. $c = 1.0$, $q_a = 1.5$.

Figure 4.1: Profit versus price $\theta$. $\sigma = 2$, $b = 1$, $s = 0.4$.

the sign of which is clearly the same as that of $-\frac{c(\sigma+1) \cdot q_a^\sigma \cdot (b+s)}{s^2 T} + \frac{\sigma}{q_m^*(0)} + [q_m^*(0)]^\sigma$. Due to the strict concavity of (4.21) in $q_m^*$, if $-\frac{c(\sigma+1) \cdot q_a^\sigma \cdot (b+s)}{s^2 T} + \frac{\sigma}{q_m^*(0)} + [q_m^*(0)]^\sigma < 0$, then the root of $-\frac{c(\sigma+1) \cdot q_a^\sigma \cdot (b+s)}{s^2 T} + \frac{\sigma}{q_m^*(\theta)} + [q_m^*(\theta)]^\sigma = 0$ must be less than $q_m^*(0)$. Thus, by the monotonically and strictly decreasing property of $q_m^*(\theta)$ in $\theta \in [\frac{c \cdot q_a^\sigma}{T} - s, b]$, the optimal $\theta^*$ that maximizes the intermediary's profit must be greater than 0. Similarly, the conditions for $\theta^* = 0$ and $\theta^* < 0$ can also be shown. ∎

Proposition 8 rigorously characterizes the conditions for subsidizing or taxing. Meanwhile, it also enables us to obtain some qualitative results regarding whether taxing or subsidizing is the optimal choice to maximize the intermediary's profit. Specifically, subsidizing should be selected if one of the following cases is satisfied:

1. Total content view $T$ (i.e., market size) is sufficiently small;

2. Production cost $c$ is sufficiently large;

3. Social benefit per content view $s$ is sufficiently small;

4. Aggregate content quality on the other platforms $q_a$ is sufficiently large;

5. Advertising revenue per content view $b$ is sufficiently large.

In the first four cases, few content producers can receive a non-negative utility by producing content without being subsidized by the intermediary (e.g., if the pro-

duction cost $c$ is very high, then content producers need to receive subsidy from the intermediary to cover part of the production cost). As a result, the intermediary cannot attract enough content views or maximize its profit without subsidizing the content producers. The last case indicates that if the intermediary can derive a sufficiently high advertising revenue per content view, then it can share the advertising revenue with the content producers to encourage more content production. Numerical results illustrating the impacts of $q_a$, $c$ and $T$ are shown in Fig. 4.1. It can be seen that the proposed payment scheme can significantly increase the intermediary's profit compared to setting $\theta = 0$ (i.e., without the payment scheme). For example, we observe from Fig. 4.1(b) that by optimally choosing the payment, the intermediary's profit increases from approximately 0.21 to 0.5 (i.e., nearly 150% increase).

Finally, we conclude this part by discussing two extreme cases, $q_a \to 0$ and $\sigma \to \infty$. When $q_a \to 0$, the *aggregate* content quality on the other platforms is negligible (e.g., very low quality or little content available). In other words, the intermediary becomes a monopolist in the market, and almost all the content views are devoted to content on the intermediary's platform. Therefore, the intermediary can tax the content producers by choosing $\theta^* \to -s$ and its profit can be arbitrarily close to $(b+s)T$. When $\sigma \to \infty$, the content becomes perfectly substitutable. Naturally, all the content views will be attracted by the content with the highest quality. This can also be verified by taking the limit $\sigma \to \infty$ in (4.6). Therefore, if $q_a > 1$ and $\sigma \to \infty$,[11] then the content produced on the intermediary's platform will receive no content views and the intermediary cannot possibly obtain a positive profit by varying $\theta$. On the other hand, if $q_a < 1$ (which is equivalent to $q_h < 1$ when $\sigma \to \infty$), then the content with a quality of 1 can receive almost all the content views and the intermediary can set $\theta^* \to -s$

---

[11]When $\sigma \to \infty$, we see from (4.3) that $q_a = \lim_{\sigma\to\infty} n_a^{\frac{1}{\sigma}} \cdot \lim_{\sigma\to\infty} q_h^{\frac{\sigma+1}{\sigma}} \cdot \lim_{\sigma\to\infty} \left[1 - \left(\frac{q_l}{q_h}\right)^{\sigma+1}\right]^{\frac{1}{\sigma}} \cdot \lim_{\sigma\to\infty} \left(\frac{1}{1+\sigma}\right)^{\frac{1}{\sigma}} = q_h$. Thus, when $\sigma \to \infty$, $q_a > 1$ if and only if the highest content quality $q_h > 1$.

to make its profit arbitrarily close to $(b + s)T$. To sum up, when $q_a \to 0$ or $\sigma \to \infty$ with $q_a < 1$, the intermediary can almost fully extract two sources of revenues, i.e., advertising and payment from content producers.

## 4.5 Extension to Heterogeneous Production Costs

In the preceding analysis, it was assumed that all the content producers incur a homogeneous production cost. We relax this assumption and generalize the preceding analysis by considering heterogeneous production costs.

To keep the analysis tractable, we assume that there are $K \geq 1$ possible values for content production costs, denoted by $c_1, c_2, \ldots, c_K$, where $0 < c_1 \leq c_2 \cdots \leq c_K$, and refer to content producers with the production cost of $c_k$ as type-$k$ content producers. Under the continuum model, the (normalized) mass of type-$k$ content producers is $n_k > 0$ such that $\sum_{k=1}^{K} n_k = 1$. To model the content quality heterogeneity, we assume that the content quality of type-$k$ content producers follows a continuous and positive PDF denoted by $f_k(q) > 0$ for $q \in [0, 1]$, and the corresponding CDF is $F_k(q)$ for $q \in [0, 1]$. Thus, the mass of type-$k$ content producers whose content quality is less than or equal to $q \in [0, 1]$ is given by $n_k F_k(q)$. As in the case of homogeneous production cost, for type-$k$ content producers, there exists marginal content quality, denoted by $q_{m_k} \in [0, 1]$, and a type-$k$ content producer with content quality greater (less) than $q_{m_k}$ will choose (not) to produce content on the intermediary's platform. For notational convenience, we use the vector expression $\mathbf{q_m} = [q_{m_1}, q_{m_2} \cdots q_{m_K}]$ wherever applicable.

### 4.5.1 Optimal Content Viewing

With heterogeneous production costs, we define the strictly concave utility function for the representative content viewer as $U(x(q), x_a \,|\, \mathbf{q_m}, q_a)$, where $x(q)$ denotes the

content view for each individual content $q \in [0, 1]$. The four properties specified in Chapter 4.3 can be similarly restated for the utility function $U(x(q), x_a \mid \mathbf{q_m}, q_a)$, and are omitted in the chapter for brevity. The quality-adjusted Dixit-Stiglitz utility function in (4.5) becomes

$$U(x(q), x_a \mid \mathbf{q_m}, q_a) = \left[ \sum_{k=1}^{K} \int_{q_{m_k}}^{1} n_k \cdot q \cdot x(q)^{\frac{\sigma-1}{\sigma}} dF_k(q) + q_a \cdot x_a^{\frac{\sigma-1}{\sigma}} \right]^{\frac{\sigma}{\sigma-1}}. \quad (4.26)$$

By maximizing (4.26) subject to $\sum_{k=1}^{K} \int_{q_{m_k}}^{1} n_k \cdot x(q) dF_k(q) + x_a \leq T$, we can derive the optimal value of $x(q)$ as

$$x^*(q) = \frac{T q^{\sigma}}{q_a^{\sigma} + \sum_{k=1}^{K} \int_{q_{m_k}}^{1} n_k \cdot q^{\sigma} dF_k(q)}, \quad (4.27)$$

for $q \in [\min\{q_{m_1}, q_{m_2}, \cdots q_{m_K}\}, 1]$, and $x^*(q) = 0$ otherwise.

### 4.5.2  Equilibrium Content Production

Following the analysis in Chapter 4.4 for homogeneous production cost, we first formally define the equilibrium marginal content quality, denoted by $\mathbf{q_m^*}$, as follows.

*Definition 2:* $\mathbf{q_m^*}$ is an *equilibrium* marginal content quality if it satisfies $\mathbf{q_m^*} = \mathbf{Q}(\mathbf{q_m^*})$, where $\mathbf{Q}(\mathbf{q_m^*}) \triangleq [Q_1(\mathbf{q_m^*}), Q_2(\mathbf{q_m^*}) \cdots Q_K(\mathbf{q_m^*})]$ is given by

$$Q_k(\mathbf{q_m^*}) = \arg\min_{q} \left\{ q \in [0, 1] : \ x^*(q \mid \mathbf{q_m^*}, q_a) \cdot (\theta + s) - c_k \geq 0 \right\}, \quad (4.28)$$

for $k = 1, 2 \cdots K$.

It can be easily shown that $\mathbf{Q}(\mathbf{q_m^*})$ is continuous in the compact convex set $\mathbf{q_m^*} \in [0, 1]^K$. Then, by Brouwer fixed point theorem [69], we know that there exists a fixed point (i.e., equilibrium point) that satisfies $\mathbf{q_m^*} = \mathbf{Q}(\mathbf{q_m^*})$. In game-theoretic analysis, uniqueness of the equilibrium point is important, as it ensures that the game has

a unique outcome. In the context of content platforms, uniqueness of the equilibrium marginal content quality allows the intermediary to maximize its long-term profit *deterministically*, since there is a unique outcome at the content production stage in response to the intermediary's payment. We first show two properties satisfied by the equilibrium point in the content production stage.

**Lemma 2.** *The equilibrium marginal content quality satisfies* $0 < q_{m_1}^* \leq q_{m_2}^* \cdots \leq q_{m_K}^* \leq 1$.

**Lemma 3.** *Let* $k^* = \max\{k = 1, 2 \cdots K \mid q_{m_k}^* < 1\}$. *The following relation is satisfied at equilibrium:*

$$\frac{x^*(q_{m_i})}{x^*(q_{m_j})} = \frac{c_i}{c_j}, \tag{4.29}$$

*for* $i, j \in \{1, 2, \cdots, k^*\}$.

Based on Lemmas 2 and 3, we obtain the following theorem regarding the existence and uniqueness of the equilibrium point $\mathbf{q_m^*} \in (0, 1]^K$.

**Theorem 5.** *For any* $\theta \in [-s, b]$, *there exists a unique equilibrium* $\mathbf{q_m^*} \in (0, 1]^K$ *in the production decision stage. Moreover,* $\mathbf{q_m^*}$ *satisfies*

$$\begin{cases} q_{m_k}^* = 1, & \text{if } x^*(1 \mid \bar{\mathbf{q}}_{\mathbf{mk}}, q_a) \cdot (\theta + s) \leq c_k, \\ q_{m_k}^* \in (0, 1), & \text{otherwise,} \end{cases} \tag{4.30}$$

*for* $k = 1, 2 \cdots K$, *where* $x^*(1 \mid \bar{\mathbf{q}}_{\mathbf{mk}}, q_a)$ *is obtained by maximizing* $U(x(q), x_a \mid \mathbf{q_m}, q_a)$ *subject to* $\sum_{k=1}^{K} \int_{q_{m_k}}^{1} n_k \cdot x(q) dF_k(q) + x_a \leq T$ *and* $\bar{\mathbf{q}}_{\mathbf{mk}} = \left[q_{m1}, q_{m2}, \cdots, q_{m(k-1)}, 1, 1, \cdots, 1\right]$ *satisfies* (4.29).

*Proof.* We first assume that there exists at least one equilibrium point $\mathbf{q_m^*} \in (0, 1]^K$. Next, we prove $q_{m_k}^* = 1$ if $x^*(1 \mid \bar{\mathbf{q}}_{\mathbf{mk}}, q_a) \cdot (\theta + s) \leq c_k$ and $q_{m_k}^* \in (0, 1)$ otherwise, for $k = 1, 2 \cdots K$, where $x^*(1 \mid \bar{\mathbf{q}}_{\mathbf{mk}}, q_a)$ is obtained by maximizing the utility function

$U(x(q), x_a | \bar{\mathbf{q}}_{\mathbf{mk}}, q_a)$ and $\bar{\mathbf{q}}_{\mathbf{mk}} = \left[ q_{m1}, q_{m2}, \cdots, q_{m(k-1)}, 1, 1, \cdots, 1 \right]$ satisfies (4.29). Since content producers with a content quality of 0 always receive negative utility if they produce, we must have $q^*_{m_k} > 0$ at equilibrium, for $k = 1, 2 \cdots K$. Suppose that $q^*_{m_k} \in (0, 1)$ (i.e., some type-$k$ content producers choose to produce content on the intermediary's platform at the equilibrium) when $x^*(1 \,|\, \bar{\mathbf{q}}_{\mathbf{mk}}, q_a) \cdot (\theta + s) \leq c_k$. Then, by Lemmas 2 and 3, type-1, type-2, $\cdots$, and type-$(k-1)$ content producers whose content qualities are sufficiently high will also choose to produce content at equilibrium. In particular, we have $\frac{x^*(q_{m_i})}{x^*(q_{m_k})} = \frac{c_i}{c_k}$, for $i = 1, 2 \cdots k - 1$. Thus, due to the "crowding effects" (Property 3), we can establish the following inequalities

$$x^*(q^*_{m_k} \,|\, \mathbf{q}^*_{\mathbf{m}}, q_a) \cdot (\theta + s) - c_k \leq x^*(1 \,|\, \bar{\mathbf{q}}_{\mathbf{mk}}, q_a) \cdot (\theta + s) - c_k \leq 0, \qquad (4.31)$$

which violate the definition of an equilibrium point in (4.28). As a result, $q^*_{m_k}$ must be 1 if $x^*(1 \,|\, \bar{\mathbf{q}}_{\mathbf{mk}}, q_a) \cdot (\theta + s) \leq c_k$. Similarly, we can prove by contradiction that $q^*_{m_k}$ must be strictly less than 1 if $x^*(1 \,|\, \bar{\mathbf{q}}_{\mathbf{mk}}, q_a) \cdot (\theta + s) > c_k$.

Now, to complete the proof of Theorem 5, it remains to show that there exists a unique equilibrium point $\mathbf{q}^*_{\mathbf{m}} \in (0, 1]^K$. Let $k^* = \max\{k = 1, 2 \cdots K \,|\, q^*_{m_k} < 1\}$, i.e., none of type-$k^*$, type-$(k^* + 1)$, $\cdots$, and type-$K$ content producers choose to produce content at equilibrium, and we can disregard these content producers without affecting the proof. Based on the monotonicity of $xq$ and Lemma 3, we can express $q^*_{m_i} = \rho(\frac{c_i}{c_{k^*}})q^*_{m_k^*}$, such that $\frac{x^*(q_{m_i})}{x^*(q_{m_k})} = \frac{c_i}{c_{k^*}}$ for $i = 1, 2 \cdots k^*$, and $\mathbf{q}^*_{\mathbf{m}} = \left[ \rho(\frac{c_1}{c_{k^*}})q^*_{m_k^*}, \rho(\frac{c_2}{c_{k^*}})q^*_{m_k^*}, \cdots, \rho(\frac{c_{k^*-1}}{c_{k^*}})q^*_{m_k^*}, q^*_{m_k^*}, 1, 1, \cdots 1 \right]$. Therefore, proving the uniqueness of $\mathbf{q}^*_{\mathbf{m}} \in (0, 1]^K$ is equivalent to proving $q^*_{m_k^*} = Q_k(\mathbf{q}^*_{\mathbf{m}})$ has a unique fixed point $q^*_{m_k^*} \in (0, 1)$. Following the same approach used in the proof of Theorem 3, we can show that $q^*_{m_k^*} = Q_k(\mathbf{q}^*_{\mathbf{m}})$ always has a unique fixed point. This proves Theorem 5. ∎

Because of the crowding effects (i.e., negative network externalities among the content producers, as specified in Property 3 in Chapter 4.3), we can easily show that

114

$x^*(1 \,|\, \bar{\mathbf{q}}^*_{\mathbf{mi}}, q_a) \geq x^*(1 \,|\, \bar{\mathbf{q}}^*_{\mathbf{mj}}, q_a)$, for $1 \leq i \leq j \leq K$. We can also see from (4.30) in Theorem 5 that type-$k$ content producers will choose to produce content at the equilibrium if and only if $\theta$ is sufficiently large such that $x^*(1 \,|\, \bar{\mathbf{q}}_{\mathbf{mk}}, q_a) \cdot (\theta + s) > c_k$, for $k = 1, 2 \cdots K$. For notational convenience, we define

$$\Theta = [\theta_0, \theta_1, \cdots \theta_K, \theta_{K+1}], \tag{4.32}$$

where $-s = \theta_0 \leq \theta_1 \leq \theta_2 \cdots \leq \theta_K \leq \theta_{K+1} = b$ and $\theta_k = \frac{c_k}{x^*(1 \,|\, \bar{\mathbf{q}}_{\mathbf{mk}}, q_a)} - s$, for $k = 1, 2 \cdots K$.[12] Thus, based on (4.30), it can be shown that if the payment set by the intermediary satisfies $\theta \in (\theta_k, \theta_{k+1}]$, then type-1, type-2 $\cdots$ and type-$k$ content producers with high content qualities will choose to produce content at equilibrium, whereas none of type-$(k + 1)$, type-$(k + 2)$, $\cdots$, and type-$K$ content producers will produce content, for $k = 1, 2 \cdots K$. As will be shown later, $\Theta$ defined in (4.32) is instrumental when we derive the intermediary's optimal payment to maximize its profit.

As in Chapter 4.4, by considering the best-response decision model, we can show that $\mathbf{q_{m,t}} = [q_{m_1,t}, q_{m_2,t} \cdots q_{m_K,t}]$, for $t = 0, 1, 2 \cdots$ evolves as follows $\mathbf{q_{m,t+1}} = \mathbf{Q}(\mathbf{q_{m,t}})$, where $\mathbf{Q}(\cdot)$ is defined in (4.28). More specific results regarding the existence and convergence of the equilibrium marginal content quality can be obtained, if we substitute the utility function $U(x(q), x_a \,|\, \mathbf{q_m}, q_a)$ with the Dixit-Stiglitz function defined in (4.26). The details are omitted for brevity.

---

[12]Subject to the value of $q_a$, we should note that $\theta_k = \frac{c_k}{x^*(1 \,|\, \bar{\mathbf{q}}_{\mathbf{mk}}, q_a)} - s$ may be greater than $b$. In such cases, for any $\theta \in [-s, b]$, type-$k$ content producers will never choose to produce content at the equilibrium. Without affecting the analysis, we can remove type-$k$ content producers from our model and consider a smaller subset of content producers. Therefore, we assume without loss of generality that the inequalities $-s = \theta_0 \leq \theta_1 \leq \theta_2 \cdots \leq \theta_K \leq \theta_{K+1} = b$ hold in the analysis.

$$\max_{q^*_{m_k} \in [q^*_{m_k}(\theta_{k+1}),1]} \left( b + s - \frac{c_k \left[ q_a^\sigma + \sum_{j=1}^k n_j \frac{1-(q^*_{m_j})^{\sigma+1}}{\sigma+1} \right]}{T(q^*_{m_k})^\sigma} \right) \cdot \bar{x} \qquad (4.33)$$

### 4.5.3 Optimal Price

As in the case of homogeneous production cost, it is cumbersome to maximize the intermediary's long-term profit without resorting to numerical results for the general utility function $U(x(q), x_a \mid \mathbf{q_m}, q_a)$. In this part, we consider the quality-adjusted Dixit-Stiglitz utility function defined in (4.26) and uniform distributions of content quality for the content producers, i.e., $f_k(q) = 1$ for $q \in [0,1]$ and $k = 1, 2 \cdots K$, and develop an algorithm to find the optimal payment in the following analysis. We first restrict the payment to $[\theta_k, \theta_{k+1}]$ such that only type-1, type-2, $\cdots$, and type-$k$ content producers[13] with high content qualities will choose to produce content at the equilibrium.

With uniform distributions of content qualities and $\theta \in [\theta_k, \theta_{k+1}]$, the optimal content view in (4.27) at the equilibrium of the content production stage becomes

$$x^*(q) = \frac{Tq^\sigma}{q_a^\sigma + \sum_{j=1}^k n_j \cdot \frac{\cdot 1-(q^*_{m_j})^{\sigma+1}}{\sigma+1}}, \qquad (4.34)$$

for $q \in [q^*_{m_1}, 1]$. Based on (4.34) and Lemma 3, we can see that $q^*_{m_j} = q^*_{m_k} \cdot \left( \frac{c_j}{c_k} \right)^{\frac{1}{\sigma}}$ at the equilibrium of the content production stage. If we express $q^*_{m_k}$ as a function of $\theta \in [\theta_k, \theta_{k+1}]$, then $q^*_{m_k}(\theta)$ is decreasing in $\theta \in [\theta_k, \theta_{k+1}]$ and $q^*_{m_k}(\theta) \in [q^*_{m_k}(\theta_{k+1}), 1]$. Since there exists no simple expression of $q^*_{m_k}(\theta)$, it is rather difficult to optimize $(b - \theta) \cdot \bar{x}$ by directly choosing the optimal $\theta^*$. Following the proof technique in Theorem 4, we can show that the profit maximization problem with heterogeneous production costs can be reformulated as (4.33), where $q^*_{m_j} = q^*_{m_k} \cdot \left( \frac{c_j}{c_k} \right)^{\frac{1}{\sigma}}$ and $\bar{x} =$

---

[13]Note that if $\theta = \theta_k$, then type-$k$ content producers with a content quality of $q = 1$ derive a zero utility and hence are indifferent between producing and not producing content at the equilibrium.

116

$T \cdot \left(1 - \frac{q_a^{\sigma}}{q_a^{\sigma} + \sum_{j=1}^{k} \frac{n_j \cdot (1 - q_{m_j}^{\sigma+1})}{\sigma+1}}\right)$. By showing the second-order derivative of (4.33) with respect to $q_{m_k}^* \in [q_{m_k}^*(\theta_{k+1}), 1]$ is strictly negative, we prove that the optimization problem in (4.33) is strictly concave in $q_{m_k}^* \in [q_{m_k}^*(\theta_{k+1}), 1]$. Thus, the unique optimal solution to (4.33) can be *efficiently* obtained. After solving (4.33), we can obtain the optimal payment as $\theta^* = \frac{c_k \left[q_a^{\sigma} + \sum_{j=1}^{k} n_j \frac{n_j 1 - (q_{m_j}^*)^{\sigma+1}}{\sigma+1}\right]}{T(q_{m_k}^*)^{\sigma}} - s$.

Next, based on the optimal solution to (4.33), we develop a recursive algorithm to find the optimal payment maximizing the intermediary's profit and describe it in Algorithm 2.

---

**Algorithm 2** Find $\theta^* \in [-s, b]$

---

$\Pi_{\mathcal{I}} \leftarrow 0$, $\theta^* \leftarrow -s$, and $k \leftarrow 1$
**while** $k \leq K$ **do**
   Solve (4.33) and denote the maximum value by $temp$
   **if** $\Pi_{\mathcal{I}} < temp$ **then**
     $\Pi_{\mathcal{I}} \leftarrow temp$
     $\theta^* = \frac{c_k \left[q_a^{\sigma} + \sum_{j=1}^{k} n_j \frac{n_j 1 - (q_{m_j}^*)^{\sigma+1}}{\sigma+1}\right]}{T(q_{m_k}^*)^{\sigma}} - s$
   **end if**
   $k++$
**end while**
**return** $\theta^*$

---

As in the case of homogeneous production cost, we can also analyze whether the intermediary should tax or subsidize the content producers. Nevertheless, we omit the result because of its similarity with Proposition 8.

## 4.6 Conclusion

In this chapter, we studied a user-generated content platform and proposed a payment scheme in which the intermediary can either tax or subsidize the content producers to maximize its profit in the presence of rational content producers and content viewers.

We first used the *representative* content viewer model to determine how the content viewers' attention is allocated across a variety of content. Then, we showed that there always exists a unique equilibrium point at which no producer can gain by changing its production decision, and that, under certain conditions, the equilibrium point is guaranteed to be reached through an iterative process in which the content producers update their decisions with limited information. Next, we formalized the intermediary's profit maximization problem and, by using the quality-adjusted Dixit-Stiglitz utility function function and the uniform distribution of content qualities as a concrete example, derived the closed-form optimal solution explicitly. We also showed the analytical conditions under which the intermediary should tax or subsidize the content producers. Then, we discussed qualitatively the impacts of the aggregate quality of content on the other platforms and content substitutability on the intermediary's profit. Finally, we generalized our model by considering heterogeneity in the content producers' production costs. Future research directions include, but are not limited to studying: (1) differentiated payment schemes in which different content producers may be taxed or subsidized differently; (2) a scenario where content producers can vary their own content quality and choose to produce on more than one content platforms; and (3) optimal payment schemes maximizing social welfare.

# CHAPTER 5

# Dynamic Scheduling and Pricing in Wireless Cloud Computing

In this chapter, we consider a wireless cloud computing system in which a profit-maximizing wireless service provider operates a data center and can provide cloud computing services to its subscribers. In particular, we focus on batch services, which, due to their non-urgent nature, allow more scheduling flexibility than their interactive counterparts. Unlike the existing research that studied *separately* demand-side management and energy cost saving techniques for the wireless cloud (both of which are critical to profit maximization), we propose a provably-efficient Dynamic Scheduling and Pricing (Dyn-SP) algorithm which, using the pricing mechanism as a lever, *proactively* adapts the service demand to workload scheduling in the data center and opportunistically utilizes low electricity prices to process batch jobs for energy cost saving. Without the necessity of predicting the future information (as assumed by some prior works), Dyn-SP can be applied to an arbitrarily random environment in which the electricity price, available renewable energy supply, wireless network capacities provided by base stations may evolve over time as an arbitrary stochastic process. It is proved that, compared to the optimal offline algorithm with future information, Dyn-SP can produce a close-to-optimal long-term profit while bounding the job queue length in the data center. We perform a simulation study based on both traces and randomly generated data to demonstrate the effectiveness of Dyn-SP. In particular, we show both analytically and numerically that a desired tradeoff between the profit and queueing
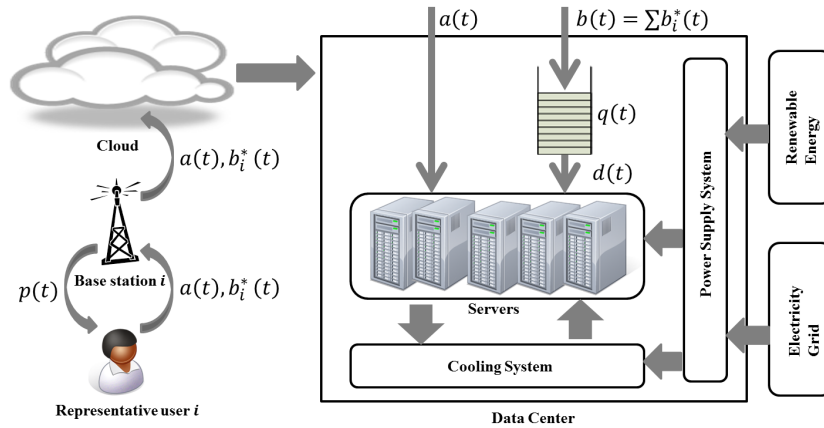
Figure 5.1: System diagram.

delay can be obtained by appropriately tuning the control parameter. Our results also indicate that, compared to the other algorithms which neglect demand-side management, cooling system energy consumption, or the queue length information, Dyn-SP achieves a higher average profit while incurring (almost) the same average queueing delay.

## 5.1 Introduction

Cloud computing has experienced an explosive growth over the past few years. As a service featuring scalability, reliability and low cost [71], "computing" is delivered to multiple clients over wireline/wireless networks, releasing clients from the hassle of maintaining their own computing infrastructure. More recently, enabled by ubiquitous communications and low-cost mobile devices, the proliferation of digital data and growing demand for outsourcing the data to a cloud for processing has attracted a significant amount of attention from various IT companies. Particularly, it has been an emerging trend that traditional wireless carriers are integrating cloud computing with wireless access services (e.g., [72]) by jointly managing their base stations and data centers.

120

In this chapter, we focus on a profit-maximizing wireless service provider that can "sell" cloud computing services to its subscribers. The system diagram is illustrated in Fig. 5.1. The wireless service provider operates multiple base stations and one data center, which has an on-site renewable energy supply system as a supplementary source of energy. Due to hardware/software constraints on mobile devices, wireless subscribers *outsource* data-intensive and/or computing-intensive applications to the cloud managed by the service provider for processing. Migrating such applications is mutually beneficial: users can access a large pool of high-performance computing resource anywhere and anytime; the wireless service provider is no longer only providing connectivity between users and the cloud service providers, but can enhance its profit by integrating cloud computing with wireless access services. Example applications that can, or even have to, be migrated to the cloud include real-time stream mining, scientific computing, visual search, and batch image processing [73]. In general, these applications can be classified as *interactive* (i.e., delay-sensitive) and *batch* (i.e., delay-insensitive). For interactive applications, the service provider needs to process them in real-time and respond immediately. By contrast, for batch applications, the service provider can defer them to an *appropriate* time instant to process, as long as the average response time constraint is satisfied. In our study, we concentrate on dynamically scheduling and pricing batch services which, due to their non-urgent nature, give the service provider more scheduling flexibility.

As a commercial entity, the wireless service provider's ultimate goal is maximizing its long-term profit subject to a service level agreement (e.g., average delay requirement) by optimally designing the demand-side management and data center operation. Nevertheless, the optimal design is challenged and complicated by the highly-intermittent nature of availability of renewable energies as well as the time-varying network capacities provided by base stations. To tackle the randomness in the environment, we propose a provably-efficient online algorithm, Dynamic Scheduling and

121

Pricing (Dyn-SP), which can be implemented using the currently available information without the necessity of predicting the future. For an arbitrarily random environment, we prove that Dyn-SP is efficient in the sense that the profit gap between the Dyn-SP and the optimal offline algorithm with $T$-slot future information is upper bounded. Meanwhile, the job queue length is also upper bounded, resulting in an finite (and bounded) average queueing delay. Moreover, the tradeoff between the queueing delay and profit can be flexibly adjusted, making Dyn-SP an appealing candidate for future wireless cloud management. We conduct extensive simulations to validate Dyn-SP. In particular, we show that the long-term profit achieved by Dyn-SP can be pushed arbitrarily close to the maximum (at the expense of increasing the queueing delay). Compared to other algorithms which neglect demand-side management, cooling system energy consumption, or the queue length information, Dyn-SP produces a higher profit while incurring the same queueing delay. Our results highlight the importance of dynamic pricing as well as the integrated approach to managing data centers (which considers both cooling system and server energy consumption). The key ideas of Dyn-SP are two-fold: **(1) reshaping the service demand:** pricing is used as a lever to proactively adapt the service demand to the time-varying data center management; and **(2) opportunistically utilizing low electricity prices:** batch jobs are deferred and processed (using electricity energy) only when the electricity price is sufficiently low relative to the length of the job queue maintained in the data center.

The rest of this chapter is organized as follows. Related work is reviewed in Chapter 5.2. Chapter 5.3 describes the model. In Chapter 5.4, we develop a provably-efficient online algorithm to maximize the service provider's average profit subject to queueing delay constraint. Simulation results are shown in Chapter 5.5 and finally, concluding remarks are offered in Chapter 5.6.

## 5.2 Related Works

In this part, we first review the existing works related to wireless spectrum management with add-value services, energy-efficient data center operation as well as pricing, and then we differentiate our study from the existing research in terms of the model and the methodology.

• Wireless spectrum management with add-value services: Wireless service providers have invested billions of dollars in the purchase of scarce radio spectrum (often through auction) from the respective authorities (e.g., Federal Communications Commission in the United States). To revert the huge investment, wireless service providers are attempting to provide new add-value services over its spectrum to the subscribers (see [74] for a quick review). For example, the advancement of cognitive radio technology makes it possible that the wireless service provider obtains additional revenues by leasing its unused spectrum to third-party users on a short-term basis [75]. In addition to spectrum leasing, relay stations can also bring additional revenues to wireless service providers: helping distant transmitters forward signals to their respective receivers at the expense of energy consumption and additional spectrum usage [76][77]. Recently, with a plethora of data and mobile applications, providing cloud-based services to wireless subscribers has attracted special attention from various wireless carriers (e.g., [72]), enabling new business opportunities. For instance, [78] proposes a new architecture migrating data-intensive and/or computation-intensive applications from mobile devices to the cloud, whereas [79] advocates the migration of energy-demanding applications from energy-constrained mobile devices to the cloud.

• Energy-efficient data center operation: Many megawatts are required to power an Internet-scale data center, invoking a growing interest in cutting electricity bills for large data centers. Several classic approaches to energy saving in data centers include power-aware architecture design and workload consolidation (e.g., [81]), cluster-

123

level (thousands of servers) power management (e.g., [82]), optimal power allocation among servers to dynamically tune the server performance (e.g., [83]), and dynamic and automated cooling system and power supply management (e.g., [84][85][86]). Another important line of research on energy-efficient data centers deals with workload management/scheduling, which is the focus of our study. For instance, [87] proposes to dynamically defer batch jobs based using a linear programming approach. [88] proposes an online right-sizing algorithm which dynamically turns on/off servers to minimizes the delay plus energy cost, under the assumption that the electricity price is fixed over time. [89] considers a similar problem but proposes to predict the future service demand using a Markov chain to determine the number of active servers. [90] quantifies the economic gains by scheduling workloads across multiple data centers, which is an empirical study without providing analytical performance bounds on the proposed scheduling algorithm. Several studies explore the opportunity of energy saving by executing jobs when and/or where the electricity prices are low (e.g., [91][92][93][94][95]). Among them, some perform local and myopic optimization at each time period without offering performance guarantees for the average energy cost or queueing delay over a large time horizon [91][92]. Other prior studies assume that the electricity price variations and/or job arrivals follow certain stationary (although possibly unknown) distributions [93][94][95]. More recently, an *integrated* approach to data center management, which takes into account server and cooling system energy consumption as well as renewable energy supply, is proposed to increase the service provider's profit [96].

• Pricing: Originated from economics literature, pricing has been studied extensively by various communities and for various purposes. In particular, pricing has been widely used as a lever to enable efficient resource management in wireless networks (e.g., [76][77]), congestion control (e.g., [98][99]), and demand coordination in smart grid [100]. In the contexts of cloud computing, [101] proposes a dynamic pric-

Table 5.1: Literature Review

| Literature | Pricing | Renewable Energy | Cooling System | Methodology | Environment | Performance Guarantee |
|---|---|---|---|---|---|---|
| [75]–[77] | Yes | No | No | Game theory | Static | No |
| [87] | No | No | No | Linear programming | Arbitrary | Yes |
| [88] | No | No | No | Greedy | Arbitrary | No |
| [89] | No | No | No | Optimization w/ future information | Arbitrary | No |
| [91] | No | Yes | No | Convex optimization | Static | No |
| [92] | No | No | No | Optimization | Static | No |
| [93]–[95] | No | No | No | Stochastic control | i.i.d. or Markovian | Yes |
| [96] | No | Yes | Yes | Optimization w/ future information | Arbitrary | No |
| **Our study** | **Yes** | **Yes** | **Yes** | **Stochastic control** | **Arbitrary** | **Yes** |

ing algorithm to regulate admission and resource allocation for social welfare maximization, [102] develops a computationally efficient and truthful auction-style pricing mechanism enabling fair competition for cloud resource among self-interested users, and [103] studies a pricing game among multiple cloud service providers and developed an algorithm to determine the minimum amount of resource to guarantee the prescribed Quality-of-Service (QoS). Combined with manufacture management (e.g., capacity sizing, workflow scheduling), pricing has also invoked lots of interest in operational research. For instance, [104] studies the equilibrium behavior of multiple users sharing a common pool of resource in the context of pricing and capacity sizing under revenue and social optimization objectives. [105] presents a decision model that integrates pricing and production decisions for the cases where the manufacturer charges the same price or different prices to different customers. [106] considers the coordination of pricing and scheduling decisions in a make-to-order environment, and proposes computationally-efficient algorithms to maximize the total net profit.

Next, we state the key differences between our study and the existing research, which are also summarized in Table 5.1.

- **First,** unlike the existing research that *separately* studied demand side management and data center operation while neglecting the interplay between workload scheduling and pricing, we propose a Dyn-SP which, using the pricing mechanism as a lever, *proactively* adapts the service demand to workload scheduling in the data

center and opportunistically utilizes low electricity prices to process batch jobs for energy cost saving. Joint design of scheduling and pricing has been touched upon mostly in operational research (e.g., [106]) which, however, cannot be applied in wireless cloud computing due to its neglectance of many important characteristics in wireless networks and data centers (e.g., time-varying wireless network capacities, electricity prices, renewable energy supply, as well as queueing delay constraint).

- **Second,** we adopt an *integrated* approach to data center management by taking into considering the server energy consumption, cooling system energy consumption as well as renewable energy supply, whereas the existing research (except for [96]) on workload scheduling neglects the cooling system energy consumption and/or renewable energy supply.

- **Last but not least,** what distinguishes our study from the prior works is the methodology. Specifically, the existing research on demand side management and workload scheduling relies on one of the following approaches: (1) optimization with predicted future information (e.g., [81][96][89] for data center management and [109]); (2) myopic optimization (e.g., [90]–[92]); (3) stochastic control techniques under the assumption that the environment, including electricity prices and renewable energy supplies, is independently and identically distributed or Markovian (e.g., [93]–[95]). By contrast, the proposed algorithm, Dyn-SP, builds on the sample-path version of Lyapunov optimization and applies for an arbitrarily random environment while providing a long-term performance guarantee.

## 5.3 Model

In this chapter, we focus on a profit-maximizing wireless service provider that dynamically prices its wireless cloud computing services. We consider a discrete-time model in which the time slots match the timescale at which the control decisions are

Table 5.2: List of Notations

| Notation | Description |
|---|---|
| $a(t)$ | Interactive service demand |
| $b(t)$ | Batch service demand |
| $d(t)$ | Processed batch jobs |
| $f(t)$ | Cooling system energy consumption |
| $y(t)$ | Renewable energy supply |
| $p(t)$ | Price for batch service |
| $u_i(b(t), t)$ | Time-dependent utility for representative user $i$ |
| $\phi(t)$ | Electricity price |
| $r(t)$ | Electricity cost |
| $h(t)$ | Profit |
| $q(t)$ | Batch job queue length |
| $C_i(t)$ | Network capacity of base station $i$ |

made. Next, we describe the service provider model as well as the user model, and then discuss the extension of our model. Key notations are listed in Table 5.2.

## 5.3.1 Service Provider

We consider a wireless service provider that delivers cloud computing services to its subscribers by jointly managing (cellular) base stations as well as data centers housing tens of thousands of servers.[1] In the following, we provide the modeling details of base stations and data centers.

### 5.3.1.1 Base station

The service provider owns $N$ base stations, indexed by $1, 2, \cdots, N$, respectively, each of which covers a certain area.[2] Assuming that certain wireless technologies (e.g., code division multiple access or CDMA) have been deployed, we use $C_i(t) \in [C_{\min}, C_{\max}]$

---

[1]Although we focus on cellular base stations, our analysis also applies to the scenario in which the service provider provides wireless access services via other technologies (e.g., Wi-Fi hot spots).

[2]In this chapter, we ignore the energy consumption of base stations, which is relatively less than that of an Internet-scale data center, whereas we leave the incorporation of energy consumption of base stations in our future work.

to represent the network capacity provided by the base station $i$ at time $t$. The network capacity is time-varying for various reasons such as the users' distances to the respective base station.

### 5.3.1.2  Data center

For the convenience of presentation, we consider that the service provider only owns one data center while noting that multiple data centers can be viewed as one *consolidated* (virtual) data center. Due to the enormous appetite for energy, a major component of operational costs of a data center is its electricity cost, which thereby has led to an increasing interest in decreasing the energy usage. Of the megawatts required to power a data center, a significant portion (typically 80%-85% [80]) is consumed by servers and cooling systems (which keeps the servers running at an appropriate temperature for reliability purposes). Despite being beyond the scope of this chapter, we note that the power system/infrastructure also incurs a non-negligible power loss (e.g., power transmission, conversion [96][97]). While significant progress has been made in improving the energy efficiency of servers and cooling systems separately (e.g., [81]–[85]), we adopt in this chapter the approach of integrated data center management [96] by jointly considering the server energy consumption and cooling system energy consumption.[3] In the following, we provide the models of the cooling system, servers, as well as the power supply system of the data center.

• Cooling system: Depending on the data center's geographic location, multiple approaches (e.g., outside air economizer, chiller plant) may be available to keep servers cool by exhausting the heat generated from densely-organized server racks to the outside. In this chapter, we consider chiller plant, which provides relatively stable cooling resources, as the only approach to cooling down servers. In general, there exists no

---

[3]As all the time slots have the same duration, we interchangeably use *energy* and *power* wherever applicable.

simple analytical relation between the number of active servers and the power consumption by the associated cooling system such that the temperature is held (relatively) constant. As a proxy, we assume that the cooling system incurs a power consumption that is linearly increasing with the number of active servers, i.e., $f(m) = \gamma \cdot m$ where $m \geq 0$ is the (normalized) number of active servers and $\gamma > 0$ is a certain constant depending on the chiller structure. Albeit simple, the linear relation has been found to be reasonably accurate through empirical measurement [96][86]. Without causing ambiguity, we sometimes use $f(t) = f(m(t))$ to represent the cooling system power consumption at time $t$. Note that our analysis still applies if a more general form of $f(t) = f(m(t))$ (e.g., convex function) is considered. Moreover, even though multiple cooling techniques are utilized to in the data center, the overall cooling energy consumption can still be expressed as a single function $f(t) = f(m(t))$ in terms of the number of active servers [96].

• Servers: To avoid delving into the intricate details of servers, we assume that the data center houses $W$ (homogeneous) servers, each of which has a normalized processing speed of 1 and incurs a power of 1 when active. Without losing key insights, this model has been widely applied in the existing literature pertaining to data center management (e.g., [96][88]). The service provider can dynamically *size* its data center by turning on/off (a certain number of) servers to adapt the server provisioning to the incoming workloads. For the sake of analysis, we ignore the toggling costs (e.g., wear-and-tear costs) incurred when servers are switched on/off, while noting that such costs can be dealt with using the approach of [88] if necessary.

• Power supply system: Typically, a modern data center is powered through multiple sources of energies such as grid power and renewable energy (in the form of solar, wind, etc.) [80]. Despite being sustainable, renewable energy supply is highly intermittent and time-varying subject to the source from which it is generated. For instance, while the solar energy supply often exhibits regular temporal variations, energy

generated from the wind is normally much less predictable. Fig. 5.4 shows the total available renewable energies in California on March 21, 2012 [112], which verifies highly intermittent nature of renewable energy supplies. Neglecting the energy emission/loss during the transmission and conversion processes, we denote in our study the available amount of renewable energy at time $t$ by $y(t) \in [0, y_{\max}]$. Note that the value of $y(t)$ is known (with a reasonably high confidence) at the beginning of time $t$ based on the current status of the renewable energy supply (e.g., the current weather condition and wind speed). For reliability issues, renewable energy is complemented by the conventional grid power, which can be purchased from utility companies at either a pre-defined contract rate or on-spot time-varying rate. In this chapter, we consider that the service provider participates in a real-time electricity market to purchase grid power at time-varying rate. Letting $\phi(t)$ be the real-time electricity price at time $t$, we can express the energy cost incurred by the purchase of grid power as $r(t) = r(\phi(t), e(t))$, where $e(t) \geq 0$ is the total amount of power purchased from utility companies (for processing batch jobs). Without considering the peak demand charging rate (which can be viewed as additional penalty charges when the data center power consumption exceeds a certain threshold), a particular form of $r(t)$ can be written as $r(t) = r(\phi(t), e(t)) = \phi(t)e(t)$, which we shall use extensively in the remainder of this chapter. Note that energy storage system (e.g., uninterruptible power supply or UPS) is often installed in a data center as an emergency power supply backup in the rare event of power outages. Nevertheless, such an energy storage system can typically only support the data center operation for a very short period of time, and the energy loss during the conversion process is often significant [80]. For such reasons, we isolate the energy storage system operation from our study without considering drawing power from there as an *additional* source of power supply, although incorporating it merely adds complexity and intricacy to the model.

Next, we specify the control decisions made by the service provider to maximize

its long-term profit. In general, there are two types of services provided by the service provider: interactive and batch. Examples of interactive services include real-time stream mining and visual search, and a distinguishable feature of interactive services is that they cannot be deferred due to the stringent delay deadlines associated. On the other hand, batch services, such as data backup services, scientific computing and batch image processing, can be delayed and processed depending on the schedule of the service provider, although the average response time is still critically important. Thus, the service provider has more flexibility in scheduling batch jobs (as long as the average response time requirement is satisfied), whereas it has to adapt its resource provisioning to the arrival of interactive jobs such that they are completed in real-time. Therefore, our focus is on the pricing decision and scheduling decision for batch jobs. In particular, we denote the service provider's price for its batch services at time $t$ by $p(t)$, which is held constant during the entire span of time $t$ but may change across time slots. For the sake of analysis, we quantify the jobs' service demand using the (normalized) number of servers. Specifically, we use $a(t) \in [0, a_{\max}]$ and $b(t) \in [0, b_{\max}]$ to represent the total service demands of interactive jobs and batch jobs, respectively, which are submitted by all the subscribers at time $t$. We later explain by which means $a(t)$ and $b(t)$ are determined. Under the discrete-time model, we assume that both interactive and batch jobs arrive at the service provider at the end of each time slot. Since batch jobs may be deferred to an *appropriate* time to process, a batch job queue is maintained in the data center, whose queue length (quantified in terms of total service demand) is represented by $q(t)$ at time $t$. By letting $d(t) \in [0, d_{\max}]$ be the amount of batch jobs processed at time $t$, we can show that the bath job queue length evolves according to the following dynamics

$$q(t+1) = \max[q(t) - d(t), 0] + b(t) = [q(t) - d(t)]^+ + b(t), \qquad (5.1)$$

with an initially empty queue (i.e., $q(0) = 0$). In the following study, we use the

queue length as a quantity indirectly indicating the average response time (which is largely dependent on the average queueing delay). Supposing that the electricity cost increases linearly with the power drawn from the grid and given the control decision $p(t)$ and $d(t)$, we can express the service provider's profit at time $t$ as

$$
\begin{aligned}
h(t) &= p(t)b(t) - r(\phi(t), e(t)) \\
&= p(t)b(t) - r(\phi(t), \left[ d(t) + f(d(t)) - [y(t) - a(t) - f(a(t))]^+ \right]^+) \quad (5.2) \\
&= p(t)b(t) - \phi(t) \cdot \left[ d(t) + f(d(t)) - [y(t) - a(t) - f(a(t))]^+ \right]^+,
\end{aligned}
$$

where $[y(t) - a(t) - f(a(t))]^+$ represents the amount of available renewable energy that can be used to process batch jobs, $f(d(t))$ and $f(a(t))$ are the cooling energy consumption for batch and interactive jobs, respectively, and $\phi(t)$ is the real-time electricity price charged by utility companies at time $t$. Note that we neglect in (5.2) the term representing the profit obtained by providing interactive services, since our focus is on batch services. Moreover, it is implicitly assumed in (5.2) that renewable energy is first dispatched to servers processing interactive jobs and then, if still available, may be utilized for batch jobs. Throughout the chapter, we use *environment* information to refer to the current electricity price $\phi_i(t)$, the amount of interactive jobs $a(t)$, and the available renewable energy supply $y(t)$.

*Remark 1:* In this chapter, we assume that the network connection between the data center and base stations (typically through broadband backbone networks) is not a bottleneck for the data transfer between the data center and end users. Hence, the maximum amount of jobs that can be submitted by end users is constrained solely by the wireless network capacity provided by the base stations. Moreover, we neglect the networking energy cost in data center operation, which, unlike the cooling system, is often independent of the server energy consumption and consumes relatively less power than servers and the cooling system [80].

### 5.3.2 Users

In practice, there are a large number of users subscribing to a wireless service provider, and it is quite cumbersome to model the (batch) service demand of each individual user in response to the price charged by the service provider. Nevertheless, the aggregate service demand of a group of users can be conveniently represented by the decision of a representative user. Thus, to determine the interactive service demand as well as the batch service demand, we adopt the widely-used representative agent model [107], which has been extensively applied in engineering contexts (e.g., Internet services [108] and communications market [99]) due to its facilitation of analysis. Specifically, in our study, all the users served by base station $i$ are consolidated into one representative user $i$, which then determines the *aggregate* service demand at each time slot by solving a utility maximization problem [107][99]. It is an intuitive understanding that the number of users is time-varying (e.g., there are more day users than night users). To account for this fact, we associate with each representative user $i$ a time-dependent utility function $u_i(b_i(t), t)$, where $b_i(t)$ is the representative user $i$'s batch service demand in response to the service provider's price $p(t)$ at time $t$. We use the time index $t$ to emphasize the time dependency of the utility function, and assume $u_i(0, t) = 0$ and that $u_i(b_i(t), t)$ is increasing in $b_i(t)$. In addition, $u_i(b_i(t), t)$ is typically concave in $b_i(t)$, which, however, is not required in our formulation. Note that we do not explicitly model the demand of interactive services (which in this chapter includes voice services) and instead treat it as a time-varying discrete-time process $\{a_i(0), a_i(1), a_i(2), \cdots\}$. Before processing the jobs, data exchange is required between the users and the base station (and then the submitted data will be transferred to the data center for processing via backbone networks). Nevertheless, the network capacity provided by a base station is limited in wireless communications, which in turn caps the maximum service demand. In particular, the maximum network capacity that can be used to support data exchange for batch jobs will be constrained by the

133

total network capacity minus the capacity used for interactive jobs. With a slight abuse of notation and to limit the number of free parameters, we use $C_i(t) \in [C_{\min}, C_{\max}]$ to represent the available network capacity provided by the base station $i$ at time $t$ for submitting batch jobs. For the convenience of presentation, we use $B(t) = B(b_i(t))$, which we assume is increasing in $b_i(t) \in [0, b_{i,\max}]$, to denote the required network capacity for submitting $b_i(t)$ amount of batch jobs. Then, we can mathematically formulate the (net) utility maximization problem for each representative user $i$ as follows

$$\max_{b_i(t) \in [0, b_{i,\max}]} \quad u_i(b_i(t), t) - p(t)b_i(t), \tag{5.3}$$

$$s.t. \qquad B(b_i(t)) \leq C_i(t). \tag{5.4}$$

Denote the optimal solution to the above problem (5.3)(5.4) by $b_i^*(t)$, which is clearly a function of $p(t)$. Then, the total service demand of all the users can be expressed as $b(t) = \sum_{i=1}^{N} b_i^*(t)$, which can alternatively be expressed as $b(t) = b(p(t))$ to stress that the batch job demand is (indirectly) determined by the service provider's pricing decision $p(t)$ at time $t$. Note that $b(t) = b(p(t))$ is decreasing in $p(t)$. Moreover, we assume without loss generality that $b(t) = b(p_{\max}) = 0$, i.e., with the maximum price charged by the service provider, no one is better off by using the batch service.

As a particular example, we shall consider $u_i(b_i(t), t) = \alpha_{i,t} \log(1 + b_i(t))$, where $\alpha_{i,t} \geq 0$ is a scalar indicating the representative user $i$'s time-dependent preference, when we derive more specific results in later parts. Without causing ambiguity, we also refer to $\alpha_{i,t}$ as *demand state*. The logarithm utility function has been widely studied in prior research (e.g., [99]). While there exists a rich body of research on modeling the utility function $u_i(b_i(t), t)$, the analytic expression of $B(t) = B(b_i(t))$ is more difficult to derive. For illustration purposes, we shall use extensively $B(t) = B(b_i(t)) = \lambda b_i(t)$ with $\lambda > 0$, which can be viewed as a first-order Taylor approximation of the actual $B(t)$ around the zero point.

*Remark 2:* An implicit assumption underlying the formulation (5.3)(5.4) is that the representative users are price takers without anticipating the impacts of their actions on the service provider's pricing decision. This is because a representative user is in essence a collection of many *individual* pricing-taking users, each of which is negligible (as in a continuum model) [99][108].

*Remark 3:* We can further consolidate $N$ representative users into one representative user, whose decision captures the decisions made by all the users subscribing to the service provider (rather than served by a certain base station). Under this model, the decision of the representative user will be denoted by a $N$-element vector, whose $i$-th element is the batch job service demand submitted through the $i$-th base station and is subject to the network capacity constraint of the $i$-th base station.

*Remark 4:* In order to maximize its long-term profit by dynamically setting prices for batch services, the service provider needs to obtain the information regarding the (representative) users' utility functions. Alternatively, as in economics literature, *demand function* (i.e., $b(t) = b(p(t))$ becomes a pre-requisite for the service provider to optimally determine its price. In practice, such information can be extracted from the history data or estimated online, as assumed in [99].

### 5.3.3 Extension to multiple service classes

To keep the model succinct and highlight our modeling framework that captures dynamic pricing and *integrated* data center management subject to network capacity constraints, we only present the basic model in the previous analysis, whereas we briefly discuss how the basic model is extended to better capture a real system. In particular, we emphasize QoS differentiation (in terms of average queueing delay) for multiple service classes.

Suppose that there are $K$ classes of services available to the users, and different ser-

vice classes may be associated with different QoS requirements. Then, each representative user $i$'s batch job demand becomes a vector $\mathbf{b_i}(t) = (b_{i,1}(t), b_{i,2}(t), \cdots, b_{i,K}(t))$, where $b_{i,k}(t)$ is the demand for class-$k$ service at time $t$, and its utility function $u_i(\mathbf{b_i}(t), t)$ is defined in terms of $\mathbf{b_i}(t)$. As an example, we can define the utility function as $u_i(\mathbf{b_i}(t), t) = \alpha_{i,t} \log(1 + \sum_{k=1}^{K} \beta_k b_{i,k}(t))$, which $\beta_k$ represents the relative importance of class-$k$ services. The service provider charges a price of $p_k(t)$ for class-$k$ services at time $t$. Thus, we can express the representative user $i$'s net utility as follows

$$u_i(\mathbf{b_i}(t), t) - \sum_{k=1}^{K} p_k(t) b_{i,k}(t), \tag{5.5}$$

which is jointly concave in $\mathbf{b_i}(t)$. The representative user $i$'s batch service demand at time $t$ is determined by maximizing (5.5) subject to $B_k(b_{i,k}(t)) \leq C_i(t)$, in which $B_k(b_{i,k}(t))$ is the required network capacity for demanding $b_{i,k}(t)$ class-$k$ batch services. In the data center, each service class is associated with a separate queue, whose queue length dynamics evolves similarly as in (5.1). For QoS differentiation, the queue length of class-$k$ batch jobs needs to be multiplied by a positive constant $\mu_k$, when the service provider decides the number of jobs to be processed based on the queue lengths. In particular, a larger value of $\mu_k$ indicates a higher level of QoS (i.e., shorter average response time).

## 5.4 Dynamic Pricing and Capacity Provisioning

This part first presents an offline optimal formulation of the problem of maximizing the service provider's long-term profit. Then, an online algorithm that can be implemented in practice is developed, followed by the proof of its efficiency with respect to the optimal offline algorithm with future information of up to $T$ time slots.

### 5.4.1 Offline problem formulation

In practice, wireless carriers stay in the market and operates data centers over a large time horizon, and hence maximizing the time-average profit is crucial for business interests. Let $\bar{h}$ be the time average of the service provider's profit $h(t)$, defined in (5.2), under a particular control policy implemented over a sufficiently large but finite time horizon with $t_{end}$ time slots:[4]

$$\bar{h} \triangleq \frac{1}{t_{end}} \sum_{\tau=0}^{t_{end}-1} h(\tau). \tag{5.6}$$

Similarly, we define $\bar{a} \triangleq \frac{1}{t_{end}} \sum_{\tau=0}^{t_{end}-1} a(\tau)$, $\bar{b} \triangleq \frac{1}{t_{end}} \sum_{\tau=0}^{t_{end}-1} b(\tau)$, and $\bar{d} \triangleq \frac{1}{t_{end}} \sum_{\tau=0}^{t_{end}-1} d(\tau)$. For notational simplicity, denote $\mathbf{z}(t) = (p(t), d(t))$ as the control decision made by the service provider at time $t$.

The problem of maximizing the service provider's long-term profit can be formulated as follows:

$$\max_{\mathbf{z}(t), t=0,1,2,\cdots,t_{end}-1} \bar{h} \tag{5.7}$$

$$\text{s.t.,} \quad \bar{b} \le \sum_{i \in \mathcal{D}_j} \bar{d}, \tag{5.8}$$

$$b(t) = \sum_{i=1}^{N} b_i^*(t), \forall t, \tag{5.9}$$

$$b_i^*(t) \text{ solves (5.3)(5.4)}, \forall t \tag{5.10}$$

$$d(t) + a(t) \le W, \forall t. \tag{5.11}$$

The constraint (5.8) specifies that the batch service demand needs to be accommodated over a large time scale (i.e., queue stability constraint), while the constraint (5.11) states that the server resource provided at any time cannot exceed the total available server resource (i.e., data center capacity constraint). Throughout the chapter, we as-

---

[4]As the time-scale duration of interest, $t_{end}$ is typically a very large value (e.g., in the order of thousands or even greater if each time slot corresponds to one hour) in practice.

sume that the problem (5.7)–(5.11) is *feasible* in the sense that there exits at least one control policy that satisfies the constraints (5.8)–(5.11). To solve the optimization problem (5.7)–(5.11), we need offline information (e.g., future renewable energy supplies, electricity prices, job demand) which, however, is unavailable in practice. Unlike the previous research that relies on (possibly inaccurate) prediction of the future environment information (e.g., [81][96][89] for data center management and [109]), we develop in the following analysis an online algorithm that makes pricing and job processing decisions based on the currently available information only.

### 5.4.2   Online Algorithm

Based on the sample-path version of recently developed Lyapunov optimization [110], we present in this part an online algorithm "Dyn-SP", whose performance is provably "good" compared to that of the optimal offline policy with $T$-slot lookahead information. The intuition of Dyn-SP is to trade the queueing delay for profit improvement by using the batch job queue length as a guidance for making pricing and scheduling decisions: batch service demand is reshaped using pricing as a lever to adapt to the data center management and batch jobs are processed only when the queue length becomes sufficiently large and/or electricity prices are sufficiently low. Therefore, the batch job queue dynamics introduced in (5.1), which specifies the queue length changes governed by the scheduling and pricing decisions, is instrumental for the service provider to make online decisions.

We describe Dyn-SP in Algorithm 1, which is purely online and requires only the current information (including available renewable energy supply, electricity price, interactive service demand) and queue lengths as the input. Solving (5.12) is a much simpler programming problem than directly solving (5.7)–(5.11). Even though (5.12) may not be convex in the pricing decision $p(t) \in [0, p_{\max}]$ (for example, $-p(t)b(t)$ may be non-convex in $p(t)$), the computation complexity of minimizing (5.12) is still

138

---
**Algorithm 3** Dyn-SP Algorithm
---
1: At the beginning of every time slot $t$, observe the current environment information (i.e., $\phi_i(t)$, $a(t)$, and $y(t)$) and the current queue length $q(t)$

2: Choose $p(t) \in [0, p_{\max}]$ to minimize

$$b(t) \cdot [q(t) - Vp(t)] = b(p(t)) \cdot [q(t) - Vp(t)], \qquad (5.12)$$

where $b(t) = b(p(t))$ is the demand function for batch services satisfying (5.9)

3: Choose $d(t) \in [0, d_{\max}]$ to minimize

$$V \cdot r\left(\phi(t), \left[d(t) + f(d(t)) - [y(t) - a(t) - f(a(t))]^+\right]^+\right) - q(t)d(t) \qquad (5.13)$$

where $\quad r\left(\phi(t), \left[d(t) + f(d(t)) - [y(t) - a(t) - f(a(t))]^+\right]^+\right) \quad = \quad \phi(t) \quad \cdot$ $\left[d(t) + f(d(t)) - [y(t) - a(t) - f(a(t))]^+\right]^+$ is the electricity cost

4: Update $q(t)$ according to (5.1)
---

affordable for the service provider because: (1) it only involves one decision variable; and (2) minimizing (5.12) is performed only once every time slot (which, in practice, may correspond to one hour). With $\phi(t) \cdot \left[d(t) + f(d(t)) - [y(t) - a(t) - f(a(t))]^+\right]^+$ being the electricity cost (for batch jobs), minimizing (5.13) belongs to linear programming and is rather easy to solve [113].

The parameter $V \geq 0$ is a control variable which we refer to as profit-delay parameter, and it can be tuned to different values to trade the queueing delay for the service provider's long-term profit. In particular, given a larger value of $V$, the service provider cares more about its profit, while the average queueing delay increases.

• Effect of $V$ on the service provider's pricing decision: Let us consider two extreme cases: $V \to 0$ and $V \to \infty$. When $V \to 0$, the batch jobs cannot tolerate any delays (i.e., essentially they become interactive jobs) and hence, as can be seen from (5.12), the service provider always sets $p(t) = p_{\max}$ such that no one uses its batch service. On the other hand, when $V \to \infty$, average queueing delay is not a concern and we can notice from (5.13) that the service provider always chooses its price $p(t) \in [0, p_{\max}]$ such that its profit $b(p(t))p(t)$ is maximized.

• Effect of $V$ on the service provider's scheduling decision: Next, let us explain how $V$ affects the tradeoff between the profit and queueing delay. For simplicity, we focus on the scenario in which the electricity price is given by $\phi(t) \cdot \left[ d(t) + f(d(t)) - [y(t) - a(t) - f(a(t))]^+ \right]^+$, where $\phi(t)$ is the real-time electricity price, $a(t)$ is the interactive service demand, and $y(t)$ is the available renewable energy supply (which is provided at no additional cost). It can be easily seen that in addition to processing interactive jobs (which naturally have a higher priority than batch jobs), renewable energy, if available, is always the first choice of energy supply to process batch jobs regardless of the batch job queue length. Nevertheless, whether or not to draw power from the electricity grid to process batch jobs depends on the current electricity price as well as the job queue length. Specifically, if $\phi(t) \leq \frac{q(t)}{V(1+\gamma)}$ is satisfied where $\gamma$ is the required cooling power per unit number of active servers (and also per unit of processed workload), then the service provider will try to schedule as many batch jobs as possible to process. Otherwise, the service provider will wait until the electricity price is *sufficiently* low relative to the job queue length to process batch jobs. On the one hand, given a large value of $V$, we see from $\phi(t) \leq \frac{q(t)}{V(1+\gamma)}$ that batch jobs are processed using power drawn from the electricity grid only when the electricity price is sufficiently low. That is, the service provider *opportunistically* utilizes low electricity prices for its batch services, which clearly reduces the energy cost (and hence increases the profit) while increasing the queueing delay. On the other hand, given a small value of $V$, batch jobs are still processed using power drawn from the electricity grid even though the electricity price is not *sufficiently low* (as can be seen from the condition $\phi(t) \leq \frac{q(t)}{V(1+\gamma)}$). Doing so will clearly reduce the queueing delay, whereas at the same time it increases the electricity cost (and hence reduces the profit). Therefore, we see that $V$ is an important parameter adjusting the tradeoff between the profit and queueing delay. In the next subsection, a more formal statement regarding the role of $V$ will be provided.

### 5.4.3 Performance Analysis

This subsection shows that the proposed online algorithm is provably-efficient against an optimal algorithm with $T$-slot lookahead information. We first describe the $T$-slot lookahead policy, and then analyze the performance of Dyn-SP against it. More specifically, we show that, given a profit-delay parameter $V$, our proposed Dyn-SP algorithm is $O(1/V)$-optimal with respect to average profit against the optimal $T$-slot lookahead policy, while the queue length is bounded by $O(V)$.

### 5.4.3.1 $T$-Slot Lookahead Policy

Here, we present the $T$-slot lookahead policy, which has full knowledge of the environment information in the next (up to) $T$ time slots. If $T$ is sufficiently large (e.g., in the extreme case $T = t_{end}$), the $T$-slot lookahead policy also "approximately" (or exactly if $T = t_{end}$) maximizes the average profit in (5.7).

We divide the time horizon of $t_{end}$ time slots into $R \in \mathbb{Z}^+$ frames, each of which contains $T$ time slots such that $t_{end} = RT$. In the $T$-slot lookahead algorithm, the service provider has future environment information of up to the next $T$ time slots and maximizes the profit subject to certain constraints. Specifically, the profit maximization problem over the $r$-th frame, for $r = 0, 1, \cdots, R - 1$, can be formulated as

$$\min_{\mathbf{z}(t), t=rT, rT+1, \cdots, rT+T-1} \frac{1}{T} \sum_{t=rT}^{t=rT+T-1} h(t) \tag{5.14}$$

$$\text{s.t.,} \quad \sum_{t=rT}^{t=rT+T-1} [b(t) - d(t)] \leq 0, \tag{5.15}$$

$$\text{Constraints } (5.9) - (5.11). \tag{5.16}$$

In the problem (5.14)–(5.16), we denote the maximum of $\frac{1}{T}\sum_{t=rT}^{t=rT+T-1} h(t)$ by $H_r^*$, which is achievable over the $r$-th frame considering all the actions including those

141

that are chosen with the perfect knowledge of environment information over the entire frame. Thus, the maximum profit over $R$ frames achieved by the optimal $T$-slot lookahead policy is

$$\frac{1}{R} \sum_{r=0}^{R-1} H_r^* .$$ (5.17)

We shall show that our online algorithm, Dyn-SP, can achieve a profit "close" to the value of (5.17).

### 5.4.3.2 Online Algorithm Analysis

Now, we present the performance analysis of our proposed online algorithm compared with the optimal $T$-slot lookahead policy.

Before showing the main theorem, we first present the slackness conditions, which bound the relationship between the required server resource and the maximum server availability. The slackness conditions are prerequisites of Theorem 1.[5]

*Slackness Conditions:* There exists a positive value $\delta > 0$ and a sequence of control decisions $\mathbf{z}(t) = (p(t), d(t))$ and $h_{i,j}(t)$ such that, for data center states $\mathbf{x}(t)$, $t = 0, 1, \cdots t_{end} - 1$, the following conditions are satisfied

$$b(t) \leq d(t) - \delta,$$ (5.18)

$$d(t) + a(t) \leq W - \delta.$$ (5.19)

We note that the above slackness conditions are not restrictive at all. On the one hand, the condition (5.18) is naturally satisfied by our formulation, as the service provider can always set $p(t) = p_{\max}$ such that $b(t) = 0$. On the other hand, the

---

[5]If we only assume that the problem (5.7)–(5.11) is *strongly* feasible without slackness conditions, the performance analysis of Dyn-SP algorithm remains similar while the upper bound on the queue length may grow as the time passes [110].

142

condition (5.19) ensures that the available server resource is always enough to process all the scheduled jobs with a certain *slackness*. In practice, this condition is quite mild and can be easily satisfied. In particular, the server resource in a data center is provisioned for the peak load, and thus the available server resource is (almost) always sufficient for processing workloads, i.e., (5.19) holds in practice. In the worst case where the data center is overloaded, admission control techniques for interactive services can be applied to complement our scheme. Note that an equivalent statement of the slackness conditions (5.18)(5.19) is $a(t) \leq W - 2\delta$. To see this point, we can choose a control sequence $\mathbf{z}(t) = (p_{\max}, \delta)$ such that (5.18) is satisfied. In addition, if $a(t) \leq W - 2\delta$ is satisfied, then clearly (5.19) is also satisfied, which establishes the equivalence of $a(t) \leq W - 2\delta$ and the slackness conditions (5.18)(5.19).

Next, we provide Theorem 1 to show a profit bound and queue length bound for Dyn-SP.

**Theorem 1.** Suppose that the slackness conditions (5.18)(5.19) are satisfied for some $\delta > 0$, that the environment information (i.e., $\phi(t)$, $a(t)$ and $y(t)$ is arbitrarily random, for $t = 0, 1, \cdots t_{end} - 1$, and that the queue length is initially zero. Then, the following statements hold.

**a.** The queue length are bounded. For any time slot $t = 0, 1, \cdots t_{end} - 1$, we have

$$q(t) \leq \frac{V A_3}{\delta},\tag{5.20}$$

where $V \geq 0$ and $A_3$ is a finite number defined in (5.35).

**b.** For any $T \in \mathbb{Z}^+$ and $R \in \mathbb{Z}^+$ such that $t_{end} = RT$, the profit achieved by Dyn-SP satisfies

$$\bar{h}^* \geq \frac{1}{R} \sum_{r=0}^{R-1} H_r^* - \frac{B + D(T-1)}{V},\tag{5.21}$$

where $\bar{h}^*$ is the (average) profit achieved by Dyn-SP for the problem (5.7)–(5.11), $B$ and $D$ are (finite) constants defined in the appendix and $H_r^*$ is the maximum profit in the $r$-th frame achieved by the $T$-slot lookahead policy.

*Proof.* Here, we provide the proof of Theorem 1. First, as a scalar measure of the queue length, we define the quadratic Lyapunov function as

$$L(q(t)) \triangleq \frac{1}{2}q^2(t). \tag{5.22}$$

Let $\triangle_T(t)$ be the $T$-slot Lyapunov drift yielded by some control policies over the interval $t, t+1, \cdots, t+T-1$:

$$\triangle_T(t) \triangleq L(q(t+T)) - L(q(t)). \tag{5.23}$$

Similarly, the 1-slot drift is

$$\triangle_1(t) \triangleq L(q(t+1)) - L(q(t)). \tag{5.24}$$

Then, it can be shown that the 1-slot drift satisfies

$$\triangle_1(t) \leq B + q(t)\left[b(t) - d(t)\right], \tag{5.25}$$

where $B$ is a constant satisfying, for all $t = 0, 1, \cdots, t_{end}$,

$$B \geq \frac{1}{2}\left[b^2(t) + d^2(t)\right], \tag{5.26}$$

where is finite due to the boundedness of $b(t)$ and $d(t)$.

144

**Part (a):** Based on (5.25), we can easily show that

$$\triangle_1(t) - V \cdot h(t) \leq B - V \cdot h(t) + q(t)\left[b(t) - d(t)\right]. \tag{5.27}$$

Thus, Dyn-SP actually minimizes the upper bound on the 1-slot Lyapunov drift minus a weighted profit shown on the right hand side of (5.27).

Let us choose a control action $\mathbf{z}'(t)$ satisfying the slackness conditions (5.18)(5.19). The corresponding 1-slot Lyapunov drift minus a weighted profit achieved satisfies

$$\triangle_1(t) - V \cdot h(t) \leq B - V \cdot h(t) - \delta \cdot q(t). \tag{5.28}$$

Since Dyn-SP minimizes the right hand side of (5.27), the 1-slot Lyapunov drift minus a weighted profit achieved by Dyn-SP must be less than or equal to that achieved by $\mathbf{z}'(t)$. In other words, the following inequality can be established

$$\triangle_1^*(t) \leq B - V \cdot (h^{\min} - h^{\max}) - \delta \cdot q(t), \tag{5.29}$$

where $h^{\max}$ and $h^{\min}$ are the maximum and minimum 1-slot profit,[6] respectively, and $\triangle_1^*(t)$ is the 1-slot Lyapunov drift achieved by Dyn-SP. Now, we define

$$P \triangleq B - V \cdot (h^{\min} - h^{\max}). \tag{5.30}$$

Thus, if the queue length $q(t)$ is greater than or equal to $P/\delta$, the 1-slot Lyapunov drift in (5.29) is non-positive. Moreover, we can show that the maximum value of the Lyapunov function is $[P/(\sqrt{2}\delta)]^2$ under the constraint that the queue length $q(t)$ is less than or equal to $P/\delta$. Thus, if the Lyapunov function is greater than $[P/(\sqrt{2}\delta)]^2$, the queue length $q(t)$ will be greater than $P/\delta$ and the Lyapunov function in the next time

---

[6] Both $h^{\max}$ and $h^{\min}$ are finite due to boundedness of $b(t)$ and $d(t)$.

slot will not increase, since the 1-slot Lyapunov drift is negative. Nevertheless, if the queue length $q(t)$ is less than or equal to $P/\delta$ during time $t$, we have

$$L(q(t+1)) \leq \frac{1}{2}\Big[q(t) + q^{diff}(t)\Big]^2 \leq L(\Theta(t)) + D + \frac{q^{diff}P}{\delta} \leq \Big(\frac{P}{\sqrt{2}\delta}\Big)^2 + D + \frac{q^{diff}P}{\delta},$$

(5.31)

where $q^{diff}(t)$ represents the absolute value of changes in the queue length, with maximum value being $q^{diff} = \max[b_{\max}, d_{\max}]$, and $D$ is a constant satisfying, for all $t = 0, 1, \cdots, t_{end} - 1$,

$$D \geq \frac{1}{2}q^{diff} \cdot \max\Big[b(t), d(t)\Big],$$

(5.32)

which is finite due to the boundedness of $b(t)$ and $d(t)$. Clearly, $L(q(0))$ satisfies (5.31), as the queue length $q(0)$ is initially zero. Then, by mathematical induction, we can show that, for any $t = 0, 1, \cdots, t_{end} - 1$,

$$L(q(t)) \leq \Big(\frac{P}{\sqrt{2}\delta}\Big)^2 + D + \frac{q^{diff}P}{\delta},$$

(5.33)

following which we see that all the queue lengths are bounded by

$$q(t) \leq \sqrt{\Big(\frac{P}{\delta}\Big)^2 + 2D + \frac{2q^{diff}P}{\delta}} = \frac{V\sqrt{\frac{P^2}{V^2} + \frac{2D\delta^2}{V^2} + \frac{2q^{diff}\delta P}{V^2}}}{\delta} = \frac{VA_3}{\delta},$$

(5.34)

where

$$A_3 = \sqrt{D_1 + D_2 + D_3},$$

(5.35)

146

in which

$$D_1 \triangleq \left[ \frac{B}{V} + h^{\max} - h^{\min} \right]^2,$$ (5.36)

$$D_2 \triangleq \frac{2D\delta^2}{V^2},$$ (5.37)

$$D_3 \triangleq \frac{2q^{diff}\delta}{V}\sqrt{D_1}.$$ (5.38)

This proves part (a) of Theorem 1.

**Part (b):** Based on (5.27), we can show that, for $r = 0, 1, \cdots, R - 1$, the $T$-slot drift minus weighted profit satisfies

$$\triangle_T^*(rT) - V \sum_{t=rT}^{rT+T-1} h^*(t) \leq BT - V \sum_{t=rT}^{rT+T-1} h(t) + \sum_{t=rT}^{rT+T-1} (t - rT)q^{diff}[b(t) - d(t)]$$
$$+ q(rT) \sum_{t=rT}^{rT+T-1} [b(t) - d(t)].$$

(5.39)

Then, after some simple mathematic manipulations based on (5.39), we can derive the following inequality

$$\triangle_T^*(rT) - V \sum_{t=rT}^{rT+T-1} h^*(t) \leq BT + V \sum_{t=rT}^{rT+T-1} h(t) + DT(T-1)$$
$$+ q(rT) \sum_{t=rT}^{rT+T-1} [b(t) - d(t)],$$

(5.40)

where $D$ is a finite constant satisfying (5.32). In (5.40), the left hand side is the $T$-slot Lyapunov drift minus weighted profit achieved by Dyn-SP, which explicitly minimizes the right hand side of (5.27). Note that the right hand side of (5.27) is smaller than or equal to that of (5.40). Thus, by considering the optimal $T$-slot lookahead policy on

147

the right hand side of (5.40), we obtain the following inequality

$$\triangle_T^*(rT) - V \sum_{t=rT}^{rT+T-1} h^*(t) \leq BT - VTH_r^* + DT(T-1) + q(rT) \sum_{t=rT}^{rT+T-1} [b(t) - d(t)]$$

$$\leq BT + VTG_r^* + DT(T-1),$$

(5.41)

where the second inequality follows from the constraint in (5.15) satisfied by the optimal $T$-slot lookahead policy. Therefore, by summing (5.41) over $r = 0, 1, \cdots, R-1$, and considering that all the queues are initially empty, it follows that
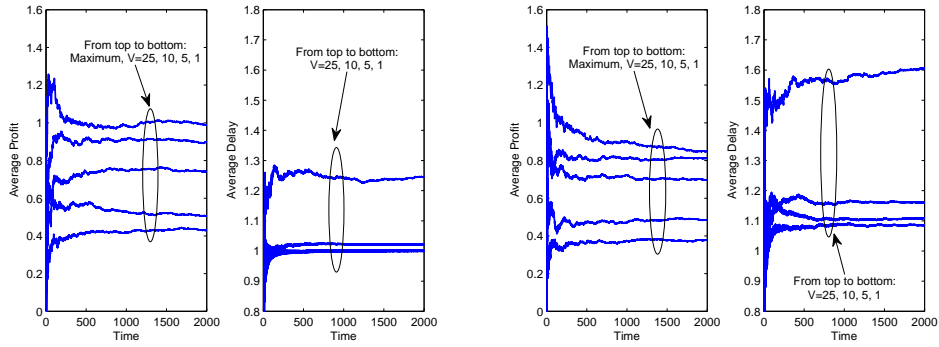
$$-V \sum_{t=0}^{RT-1} h^*(t) \leq BTR - VT \sum_{r=0}^{R-1} H_r^* + RDT(T-1) - L(q(RT-1))$$

$$\leq BTR - VT \sum_{r=0}^{R-1} H_r^* + RDT(T-1).$$

(5.42)

Finally, by dividing both sides in (5.42) by $VTR$, we have

$$\bar{h}^* = \frac{1}{TR} \sum_{t=0}^{RT-1} h^*(t) \geq \frac{1}{R} \sum_{r=0}^{R-1} H_r^* - \frac{B + D(T-1)}{V},$$

(5.43)

which shows that the online algorithm can achieve a profit within $O(1/V)$ to the minimum cost achieved by the optimal $T$-slot lookahead policy. This proves part (b) of Theorem 1. ∎

Theorem 1 shows that, given a profit-delay parameter $V$, our algorithm Dyn-SP is $O(1/V)$-optimal with respect to the average profit against the optimal $T$-slot lookahead policy, while the queue length is bounded by $O(V)$. More specifically, the inequality (5.20) bounds the queue length: the queue length (which is closely related to the average queueing delay) is bounded by $O(V)$ where $V$ is the profit-delay parameter in Algorithm 1. The queue length bound is tighter when $V$ is smaller. The inequality (5.21) shows that the average profit is bounded within an additional $O(1/V)$ profit

(a) $a(t)$ is uniformly distributed between $0$ and 3.

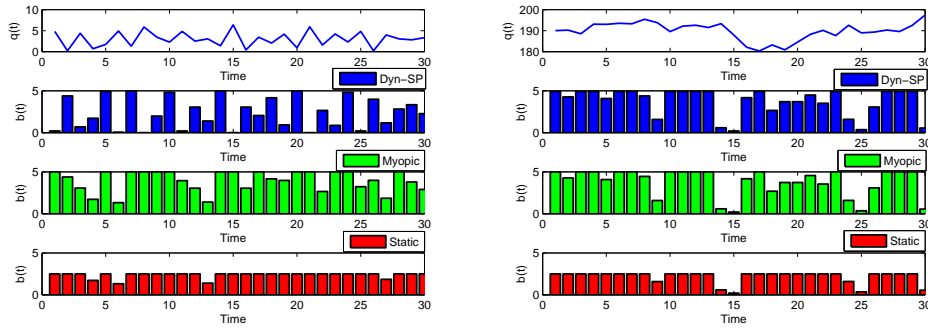(b) $a(t)$ is uniformly distributed between $0$ and 8.

Figure 5.2: Dyn-SP with different $V$.

loss of that achieved by the optimal $T$-slot lookahead policy. The profit loss is smaller when $V$ is larger. Note that the upper bound on the profit loss $\frac{B+D(T-1)}{V}$ grows with the value of $T$ (i.e., information window size of the *oracle* which owns the future information), which can be intuitively understood as follows: with more information about the future environment, the profit achieved by the offline algorithm increases, thereby resulting in an enlarged profit gap.

The key insight obtained from Theorem 1 is that, by appropriately tuning the profit-delay parameter $V$, we can achieve a desired tradeoff between the profit and queue length with analytic performance guarantees. Finally, it should be pointed out that although the derived analytic performance bound, which applies for an arbitrarily random environment subject to mild slackness conditions, may not be very *tight* for certain system settings, the actual performance of Dyn-SP is reasonably good in practice and outperforms the other algorithms as shown in the next part.

## 5.5 Numerical Results

We perform a simulation study to evaluate our algorithm using traces of hourly electricity prices and renewable energy supplies for a data center operating in California.

149

(a) Batch service demand under different pric-  (b) Batch service demand under different pric-
ing schemes with $V = 10$.                       ing schemes with $V = 1000$.

Figure 5.3: 30-hour snapshot.

We conduct three sets of experiments:

- Profit-maximization with different $V$: study the impact of the profit-delay parameter $V$ on profit maximization.

- Algorithm comparison: compare Dyn-SP with three other algorithms (i.e., Stat-Pricing, ServerOnly, BestEffort) explained later.

- Extension to multiple service classes.

The experimental results show that (1) Dyn-SP maximizes profit while bounding the queue length by opportunistically pricing and scheduling batch jobs when the electricity prices are sufficiently low relative to the queue length; (2) Compared to the other algorithms, Dyn-SP can significantly increase the service provider's profit while still providing a queueing delay guarantee; (3) Dyn-SP can be extended to multiple service classes for QoS differentiation.

### 5.5.1  Setup

In this part, we introduce *default* settings that are used throughout the simulations unless otherwise stated. We consider a time-horizon of 2000 time slots.

- Data center: We re-scale real-world traces of hourly electricity prices and renew-
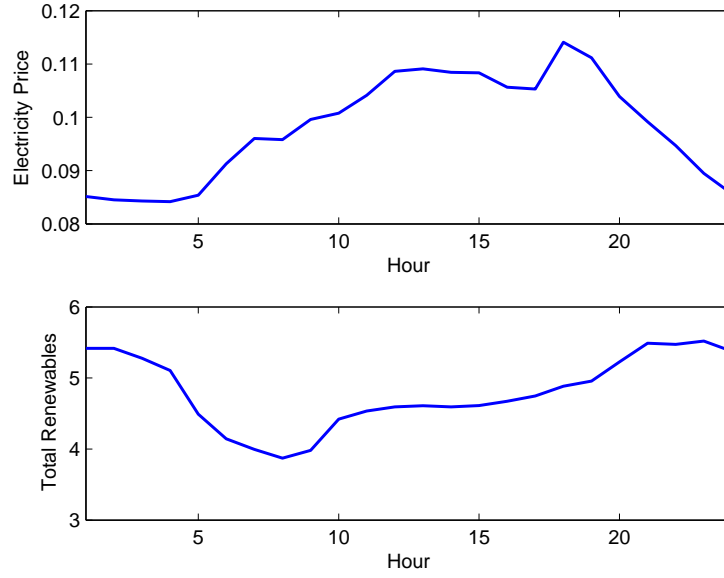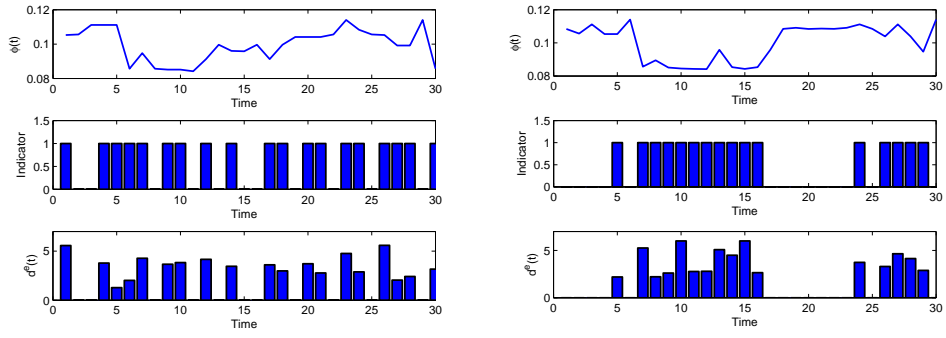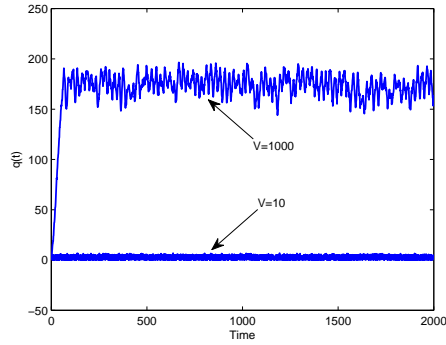
Figure 5.4: Traces of (normalized) hourly electricity prices and total available renewable energies (March 21, 2012) in California, USA [112]

able energy supplies in California, USA [112]. A 24-hour trace sample is plotted in Fig. 5.4. The total number of servers in the data center is normalized to 10. Based on empirical data provided in [80], we assume that the cooling system consumes a power equal to 0.75 times the server power, i.e., $f(m) = 0.75 \cdot m$, where $m$ is the number of active servers.

• User: For simplicity, we assume that there are 10 base stations (each of which may correspond to multiple physical base stations in practice). The utility function for the representative user $i$ is $u_i(b_i(t), t) = \alpha_{i,t} \log(1 + b_i(t))$, in which the demand state $\alpha_{i,t}$ is assumed to be independently and uniformly distributed in $[0, 1]$. The service demand for interactive applications, i.e., $a(t)$, is modeled as a uniformly distributed random variable in $[0, 8]$. We assume that a unit of network capacity is required for submitting a unit of batch service demand to the service provider, and that the available network capacity $C_i(t)$ for batch services follows a uniform distribution in $[0, 10]$.

(a) Processed batch jobs using electricity grid (b) Processed batch jobs using electricity grid
with $V = 10$.                                      with $V = 1000$.



(c) Queue length traces with $V = 10$ and $V = 1000$.

Figure 5.5: 30-hour snapshot and queue length trace.

Note that we apply the above settings to demonstrate the effectiveness of our proposed Dyn-SP, whereas our analysis also applies to any other settings.
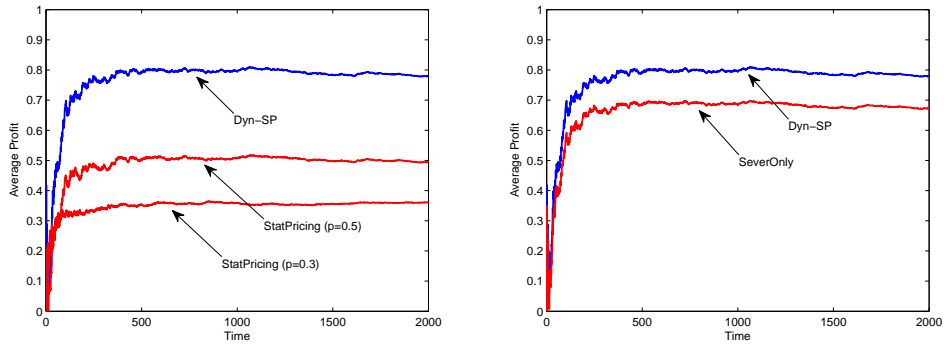
### 5.5.2 Experimental Results

In this part, we provide detailed experimental results based on the system settings described above.[7]

---

[7]The average values at time $t = 1, 2, \cdots$ are obtained by summing up all the values up to time $t$ and then dividing the sum by $t$.
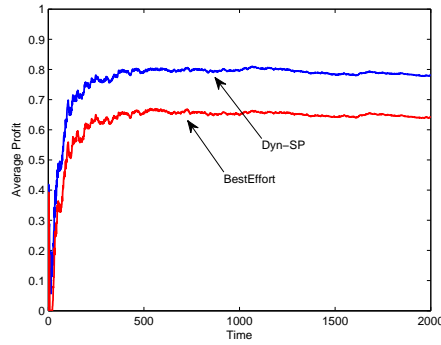
#### 5.5.2.1 Profit-maximization with different $V$

We show in Fig. 5.2 the average profit and average delay achieved by Dyn-SP with various values of $V$. It can be seen that when $V$ increases, the average profit increases at the expense of increasing the average queueing delay, which verifies Theorem 1. We note that the maximum profit can only be obtained in an online fashion if there is no queueing delay constraint. In fact, the average delay corresponding to the maximum profit in Fig. 5.2 approaches to *infinity* and thus, it is not shown. Moreover, we do not show the profit achieved by the optimal offline algorithm with $T$-slot lookahead information, which clearly cannot outperform the maximum profit. In Fig. 5.2(a) and Fig. 5.2(b), the interactive service demand $a(t)$ follows a uniform distribution in $[0, 3]$ and $[0, 8]$, respectively. Thus, in Fig. 5.2(a), there is more renewable energy available for processing batch jobs, which, as can be seen by comparing Fig. 5.2(a) and Fig. 5.2(b), reduces the long-term electricity cost.

Next, let us illustrate how Dyn-SP *reshapes* the service demand to adapt it to the data center management and how it *opportunistically* utilizes low electricity prices. First, we show in Fig. 5.3 a 30-hour snapshot of queue lengths and batch service demand under various pricing schemes (i.e., Dyn-SP, Myopic which maximizes the instantaneous profit without considering the impact of the pricing decision on the queue dynamics, and Static which uses a fixed price of $p(t) = 0.4$, for $t = 0, 1, 2, \cdots 1999$). Note that in Fig. 5.3, the demand state $\alpha_{i,t}$ takes a constant value of $0.5$ to isolate the demand state randomness and to highlight the role of queue length in setting the price. We observe from Fig. 5.3(a) that Dyn-SP can adapt the demand to the queue length: in general, the service provider will set a lower price to attract a higher batch service demand when the queue length is smaller, whereas a higher price will be used to suppress the batch service demand when the queue length is larger (such that excessive delay can be avoided). By contrast, Myopic and Static pricing schemes ignores the queue length information and thus, the batch service demands under these two pricing

(a) Dyn-SP versus StatPricing.



(b) Dyn-SP versus ServerOnly.



(c) Dyn-SP versus BestEffort.

Figure 5.6: Algorithm Comparison.

schemes are only constrained by the available network capacity (in addition to user rationality, which is captured by net-utility maximization formulated in (5.3)(5.4)). As we noted Chapter 5.4, giver a larger value of $V$, the service provide cares more about its profit and less about the queueing delay. Fig. 5.3(b) reflects this point by showing that Dyn-SP almost neglects the queue backlog and sets a price, which is nearly the same as that set by Myopic, when $V$ is sufficiently large.

Now, let us see how Dyn-SP adjusts the tradeoff between the electricity cost saving and queueing delay. Fig. 5.5 shows a 30-hour snapshot of the amount of processed batch jobs using energy drawn from the electricity grid. We have stated in Chapter 5.4 that the electricity grid supplies power for batch services if and only if $\phi(t) \leq \frac{q(t)}{V(1+\gamma)}$ is satisfied. We define an indicator which is equal to one if $\frac{q(t)}{V(1+\gamma)}$ holds. Fig. 5.5(a) ver-

ifies that batch jobs are processed using electricity energy if and only if $\frac{q(t)}{V(1+\gamma)}$ holds, i.e., the electricity price is sufficiently low relative to the queue length. In particular, we notice from Fig.5.5(a) that, in general, Dyn-SP will draw more electricity energy to process batch jobs when the electricity price is lower, i.e., low electricity prices are *opportunistically* utilized. This phenomenon is better reflected in Fig. 5.5(b), in which $V = 1000$ indicates that the service provider will shift the processing of batch jobs to time slots with low electricity prices. However, doing so will inevitably increase the queue length (as shown in Fig. 5.5(c)) as well as the queueing delay. Therefore, Fig. 5.5 illustrates the key insight of Dyn-SP: processing batch jobs using electricity energy only when the electricity price is sufficiently low relative to the queue length.
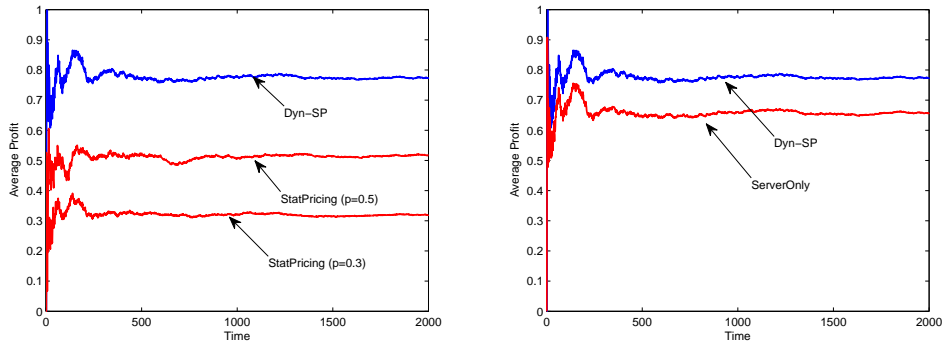
### 5.5.2.2   Algorithm comparison

We next compare Dyn-SP against three other algorithms described as follows.

- StatPricing: StatPricing uses a fixed price at all times, while it uses Step 3 in Algorithm 1 to make scheduling decisions.

- ServerOnly: ServerOnly ignores the cooling system energy consumption and only considers server energy consumption when making scheduling decisions.[8]

- BestEffort: BestEffot tries to maximize the service provider's instantaneous profit and process the submitted batch jobs as soon as possible.

We show Fig. 5.6 the profit comparison subject to (almost) the same average delay constraint of 1.15.[9] It shows that Dyn-SP outperforms all the other algorithms in terms of the average profit. Moreover, Fig. 5.6(a) shows the importance of dynamic pricing, which reshapes the batch service demand to the data center operation, while Fig. 5.6(b) indicates that an integrated approach to managing the data center, i.e., considering both cooling system and server energy consumption, is important for profit
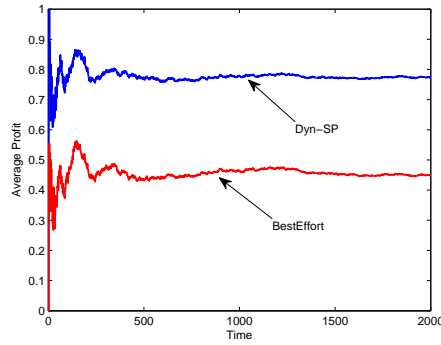
---

[8]Cooling system energy consumption is still added when the scheduling decision is made.
[9]Note that due to the discrete-time model, the minimum queuing delay is one time slot.

(a) Dyn-SP versus StatPricing.

(b) Dyn-SP versus ServerOnly.

(c) Dyn-SP versus BestEffort.

Figure 5.7: Algorithm Comparison with uniformly distributed distributed electricity price and renewable energy supply.

maximization. As shown in Fig. 5.6(c), the BestEffort algorithm is inferior to Dyn-SP in terms of profit maximization, since it neglects the queue length information and hence it cannot reshape the batch job service demand to the data center operation or exploit low electricity prices. To show the robustness of Dyn-SP against the random environment, we consider uniformly distributed electricity prices and available renewable energy supplies. In particular, the electricity price $\phi(t)$ is uniformly distributed in $[0, 0.25]$, and the available renewable energy supply $y(t)$ is uniformly distributed in $[0, 6]$. Under this setting, Fig. 5.7 shows that Dyn-SP can still outperform the other three algorithms (i.e., StatPricing, ServerOnly, and BestEffort) in terms of the average profit while incurring the same average queueing delay.
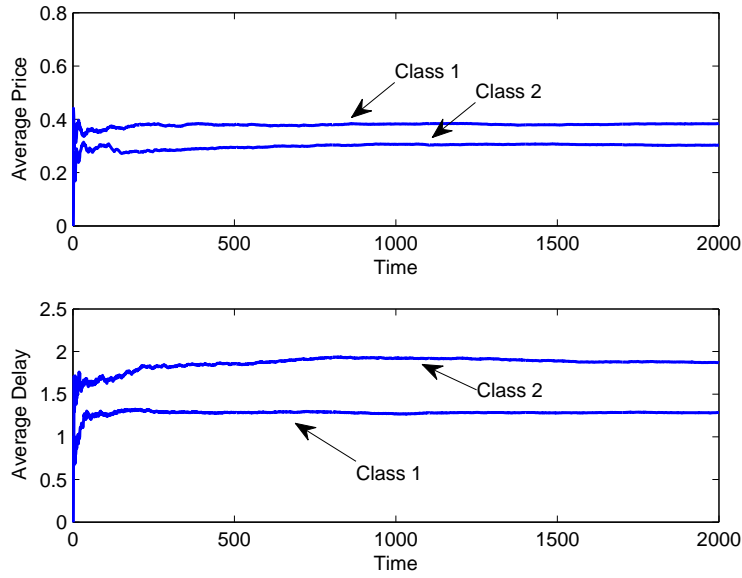
Figure 5.8: Average price and queueing delay for two service classes.

### 5.5.2.3 Extension to multiple service classes

The proposed Dyn-SP can be extended to multiple service classes and provided differentiated QoS in terms of the average queueing delay. For the ease of illustration, we consider only two service classes, and assume that the utility function of the representative user $i$ is $u_i(\mathbf{b_i}(t), t) = \alpha_{i,t} \log(1 + b_{i,1}(t) + 2b_{i,2}(t))$. Without loss of generality, we assume that class-2 service requires a lower queueing delay on average than class-1 service. Intuitively, the service provider sets a higher price for class-2 service than class-1 service. This can be seen from Fig. 5.8, in which it also shows that the average queueing delay for class-2 service is lower than that for class-1 service. Because of space limitations, we omit the results showing the tradeoff between the average profit and queueing delay, which are similar with Fig. 5.2.

157

## 5.6 Conclusion

In this chapter, we considered a profit-maximizing wireless service provider operating a data center and "selling" cloud computing services to its subscribers. The focus of our study was on batch services, which do not have urgent delay constraints. Wireless subscribers were modeled using the representative agent model, whose service demand is influenced by the service provider's pricing decision. To maximize the service provider's long-term profit, we proposed a provably-efficient online algorithm Dyn-SP which can be implemented based on the currently available information. Dyn-SP is applicable to an arbitrarily random environment in which the electricity price, available renewable energy supply, wireless network capacities provided base stations may evolve over time as an arbitrary stochastic process. We proved that, compared to the optimal offline algorithm with future information, Dyn-SP can produce a close-to-optimal long-term profit while bounding the job queue length in the data center. The key idea of Dyn-SP is: (1) using pricing to *proactively* adapt the service demand to workload scheduling in the data center; and (2) opportunistically utilizes low electricity prices to process batch jobs for energy cost saving. We performed a simulation study to demonstrate the effectiveness of Dyn-SP. In particular, it was shown both analytically and numerically that a desired tradeoff between the profit and queueing delay can be obtained by appropriately tuning the control parameter. Our results also indicated that, compared to the other algorithms which neglect demand-side management, cooling system energy consumption, or the queue length information, Dyn-SP achieves a higher average profit while incurring (almost) the same average queueing delay.

# CHAPTER 6

# Conclusion

In this dissertation, we developed a rigorous and formal framework that integrates strategic pricing with system resource management for a wide range of computer and communications systems including wireless cooperative relay networks, wireless communications markets, online user-generated content platforms, and mobile cloud computing systems. Instead of directly assuming a demand function of the price, the framework explores and explicitly considers the dynamic process in which self-interested users strategically interact with each other and respond to the charged price. We showed that the users' self-interested behaviors can be aligned with the system designer's goal via appropriate pricing mechanisms. Furthermore, in the presence of user heterogeneity, pricing was used to proactively reshape the users' behaviors/demands and adapt them to the system resource management. The framework was shown to be applicable to an arbitrarily random environment.

As the first application of the proposed framework, we considered wireless cooperative relay networks in which one dedicated relay node consumes its own limited power to forward multiple sources' signals to their respective destinations. Both uniform and differentiated pricing algorithms were proposed to maximize the system utility that can be defined in any form. The proposed pricing algorithms reimburses the relay for its power consumption and, if appropriately set, can enforce the users to transmit at desired power levels even though they are self-interested. Then, we studied wireless communications markets where the subscribers are heterogeneous in terms of their valuations of services as well as their data service demands. We proposed a gen-

eral yet practical pricing scheme that includes unlimited pricing, usage-based pricing, and capped pricing as special cases. A computationally-efficient algorithm was proposed to find the sub-optimal pricing scheme maximizing the service provider's long-term profit. Next, We turned to another important applications: online user-generated content platforms. Modeled as a two-sided market, the user-generated content platform allows users to exchange their content and also possibly make profits, as dictated by the pricing algorithm set by the intermediary who monetizes the platform. Our proposed pricing algorithm explicitly considered the participants' rationality as well as the content substitutability, which are key features of online content platforms. Finally, we focused on a wireless cloud computing system in which a profit-maximizing wireless service provider operates a data center and can provide cloud computing services to its subscribers. We proposed a novel dynamic pricing and scheduling algorithm which, using the pricing mechanism as a lever, proactively adapts the service demand to workload scheduling in the data center and opportunistically utilizes low electricity prices to process batch jobs for energy cost saving. Without the necessity of predicting the future information (as assumed by some existing works), the proposed algorithm is provably efficient in the sense that the resulting performance loss with respect to the optimal offline algorithm with perfect future information is upper bounded.

The proposed design framework can be extended in three directions. First, it can be generalized to incorporate several important issues such as price anticipation. An implicit assumption underlying our framework is that users are not price-anticipating, i.e., users are price takers. While this is a reasonable assumption in the sense that there are many users each of which only has a negligible impact on the service provider's pricing decision, investigating pricing anticipation is important as part of the sensitivity study (especially when the number of users in the system is quite limited). Second, the proposed design framework can be applied in several emerging large-scale resource management scenarios including utility-based micro-grid management and intelligent

energy distribution. For instance, recent advances in smart meters and communication technologies will allow utility companies to proactively modify the energy distribution load and demand at household levels by using price signals. Together with the increasing penetration of renewable energy sources, the ability to reshape the customer demand poses both opportunities and challenges in intelligently managing the energy distribution to satisfy the soaring energy demand while utilizing the clean energy as much as possible. The proposed framework is a promising solution providing an intellectual guidance for future smart grid designs. Third, the proposed design framework can be implemented in real systems for a practical understanding of implementation issues. Moreover, it can be unified with other relevant theories (e.g., cross-layer design in communications systems) to advance the theoretic framework of jointly designing economic incentives and resource management.

# REFERENCES

[1] B. Rankov and A. Wittneben, "Spectral efficient protocols for half-duplex fading relay channels," *IEEE J. Sel. Areas Commun.*, vol. 25, no. 2, pp. 379-389, Feb. 2007.

[2] A. Sendonaris, E. Erkip, and B. Aazhang, "User cooperative diversity - Part I: System description," *IEEE Trans. Commun.*, vol. 51, no. 11, pp. 1927-1938, Nov. 2003.

[3] J. N. Laneman, D. N. C. Tse, and G. W. Wornell, "Cooperative diversity in wireless networks: Efficient protocols and outage behavior," *IEEE Trans. Inform. Theory*, vol. 50, no. 12, pp. 3062-3080, Dec. 2004.

[4] S. Ren and K. B. Letaief, "Minimum sum expected distortion in cooperative networks," in *Proc. IEEE Intl. Conf. Commun.*, June 2009.

[5] S. Serbetli and A. Yener, "Relay assisted F/TDMA ad hoc networks: Node classification, power allocation and relaying strategies," *IEEE Trans. Commun.*, vol. 56, no. 6, pp. 937-947, June 2008.

[6] M. Chen and A. Yener, "Power allocation for F/TDMA multiuser two-way relay networks," submitted to *IEEE Trans. Wireless Commun.*, March 2009.

[7] M. Chen and A. Yener, "Multiuser two-way relaying: Detection and interference management," *IEEE Trans. Wireless Commun.*, vol. 8, no. 8, pp. 4296-4305, Aug. 2009.

[8] E. Altman, T. Boulogne, R. El-Azouzi, T. Jiménez, and L. Wynter, "A survey on networking games in telecommunications," *Computer Operations Research*, vol. 33, no. 2, pp. 286-311, Feb. 2006.

[9] A. A. Daoud, T. Alpcan, S. Agarwal, and M. Alanyali, "A stackelberg game for pricing uplink power in wide-band cognitive radio networks," in *Proc. IEEE Conf. Decision and Control*, Dec. 2008.

[10] J. Zhang and Q. Zhang, "Stackelberg game for utility-based cooperative cognitive radio networks," in *Proc. ACM MobiHoc*, May 2009.

[11] B. Wang, Z. Han, and K. J. R. Liu, "Stackelberg game for distributed resource allocation over multiuser cooperative communication networks," *IEEE Trans. Mobile Computing*, vol. 8, no. 7, pp. 975-990, July 2009.

[12] N. Shastry and R. S. Adve, "Stimulating cooperative diversity in wireless ad hoc networks through pricing," in *Proc. IEEE Intl. Conf. Commun.*, June 2006.

[13] O. Ileri, S. C. Mau, and N. B. Mandayam, "Pricing for enabling forwarding in self-configuring ad hoc networks" *IEEE J. Sel. Areas Commun.*, vol. 23, no. 1, pp. 151-162, Jan. 2005.

[14] N. Feng, S. C. Mau, and N. B. Mandayam, "Pricing and power control for joint network-centric and user-centric radio resource management," *IEEE Trans. Commun.*, vol. 52, no. 9, pp. 1547-1557, Sep. 2004.

[15] N. Feng, S. C. Mau, and N. B. Mandayam, "Joint user-centric and network-centric radio resource management in a multicell system," *IEEE Trans. Commun.*, vol. 53, no. 7, pp. 1114-1118, July 2005.

[16] J. Huang, R. Berry, and M. Honig, "Auction-based spectrum sharing," *ACM/Springer J. Mobile Networks and Applications*, vol. 11, no. 3, pp. 405-418, June 2006.

[17] J. Huang, Z. Han, M. Chiang, and H. V. Poor, "Auction-based resource allocation for cooperative communications," *IEEE J. Sel. Areas Commun.*, vol. 26, no. 7, pp. 1226-1237, Sep. 2008.

[18] C. U. Saraydar, N. B. Mandayam, and D. J. Goodman, "Efficient power control via pricing in wireless data networks," *IEEE Trans. Commun.*, vol. 50, no. 2, pp. 291-303, Feb. 2002.

[19] F. P. Kelly, "Charging and rate control for elastic traffic," *European Trans. Telecommun.*, vol. 8, no. 1, pp. 33-37, Jan.-Feb. 1997.

[20] F. P. Kelly, A. Mauloo, and D. Tan, "Rate control in communication networks: Shadow prices, proportional fairness, and stability," *J. of the Operational Research Society*, vol. 49, pp. 237-252, 1998.

[21] R. Yates, "A framework for uplink power control in cellular radio systems," *IEEE J. Sel. Areas Commun.*, vol. 13, no. 7, pp. 1341-1347, Sep. 1995.

[22] T. Holliday, A. Goldsmith, and P. Glynn, "Distributed power control for time varying wireless networks: Optimality and convergence," in *Proc. Allerton Conf. Commun., Control, and Computing*, pp. 1024-1033, Oct. 2003.

[23] W. Yu, G. Ginis, and J. Cioffi, "Distributed multiuser power control for digital subscriber lines," *IEEE J. Sel. Areas Commun.*, vol. 20, no.5, pp. 1105-1115, June 2002.

[24] W. Yu, "Multiuser water-filling in the presence of crosstalk," in *Inform. Theory and Applications Workshop*, Feb. 2007.

[25] T. Başar and R. Srikant, "Revenue-maximizing pricing and capacity expansion in a many-users regime," in *Proc. IEEE Infocom* pp. 1556-1563, June 2002.

[26] U. O. Candogan, I. Menache, A. Ozdaglar, and P. A. Parrilo, "Near-optimal power control in wireless networks: A potential game approach," in *Proc. of IEEE Infocom*, Mar. 2010.

[27] P. Marbach and R. Berry, "Downlink resource allocation and pricing for wireless networks," in *Proc. IEEE Infocom*, pp. 1470-1479, June 2002.

[28] C. Shen and M. van der Schaar, "Optimal resource allocation for multimedia applications over multiaccess fading channels," *IEEE Trans. Wireless Commun.*, vol. 7, no. 9, pp. 3546-3557, Sep. 2008.

[29] S. Ren and M. van der Schaar, "Distributed power allocation in multi-user multi-channel cellular relay networks," *IEEE Trans. Wireless Commun.*, vol. 9, no. 6, pp. 1952-1964, June 2010.

[30] M. Dohler and A. H. Aghvami, "A crash-course on cooperative wireless networks," in *IEEE Intl. Conf. Commun. Half-Day Tutorial*, May 2008.

[31] Y. Su and M. van der Schaar, "A new perspective on multi-user power control games in interference channels," *IEEE Trans. Wireless Commun.*, vol. 8, no. 6, pp. 2910-2919, June 2009.

[32] S. Bhashyam, A. Sabharwal, and A. Aazhang, "Feedback gain in multiple antenna systems," *IEEE Trans. Commun.*, vol. 50, no. 5, pp. 785-798, May 2002.

[33] H. R. Varian, *Microeconomic Analysis*. W. W. Norton & Company, 1992.

[34] M. J. Osborne and A. Rubinstein, *A Course in Game Thoeory*. MIT Press, Cambridge, MA, 1994.

[35] S. Ren, J. Park, and M. van der Schaar "User subscription dynamics and revenue maximization in communication markets" *IEEE Infocom*, Apr. 2011.

[36] A. Zemlianov and G. de Veciana, "Cooperation and decision making in wireless multiprovider setting," *IEEE Infocom*, Mar. 2005.

[37] L. He and J. Walrand, "Pricing and revenue sharing strategies for Internet service providers," *IEEE J. Sel. Areas Commun.*, vol. 24, no. 5, pp. 942-951, May 2006.

[38] M. Manshaei, J. Freudiger, M. Felegyhazi, P. Marbach, and J. P. Hubaux, "On wireless social community networks," *IEEE Infocom*, Apr. 2008.

[39] C. K. Chau, Q. Wang, and D. M. Chiu, "On the viability of Paris metro pricing for communication and service networks," *IEEE Infocom*, Mar. 2010.

[40] N. Shetty, S. Parekh, and J. Walrand, "Economics of femtocells," *IEEE Globecom*, Dec. 2009.

[41] N. Shetty, G. Schwartz, and J. Walrand, "Internet QoS and regulations," *IEEE/ACM Trans. Networking*, vol. 18, no. 6, pp. 1725-1737, Dec. 2010.

[42] G. Kesidis, A. Das, and G. de Veciana, "On flat-rate and usage-based pricing for tiered commodity internet services," CISS, 2008.

[43] J. Musacchio and D. Kim, "Network platform competition in a two-sided market: implications to the net neutrality issue," *TPRC: Conf. on Commun., Inform., and Internet Policy*, Sep. 2009.

[44] T. Basar and R. Srikant, "Revenue-maximizing pricing and capacity expansion in a many-users regime," in *Proc. IEEE Infocom*, Jun. 2002.

[45] L. Jiang, S. Parekh, and J. Walrand, "Time-dependent network pricing and bandwidth trading," *IEEE Intl. Workshop Bandwidth on Demand (BoD)*, 2008.

[46] D. Acemoglu and A. Ozdaglar, "Competition and efficiency in congested markets," *Mathematics of Operations Research*, vol. 32, no. 1, pp. 1-31, 2007.

[47] Z. Kun, D. Niyato, and P. Wang, "Optimal bandwidth allocation with dynamic service selection in heterogeneous wireless networks," to appear in *IEEE Globecom*, Dec. 2010.

[48] S. Shakkottai, E. Altman and A. Kumar, "Multihoming of users to access points in WLANs: A population game perspective," *IEEE J. Sel. Areas Commun.*, vol. 25, no. 6, pp. 1207-1215, Aug. 2007.

[49] V. Krishnamurthy, "Decentralized spectrum access amongst cognitive agents – An interacting multivariate global games approach," *IEEE Trans. Signal Processing*, vol. 57, no. 10, pp. 3999-4013, Oct. 2009.

[50] R. Gibbens, R. Mason, and R. Steinberg, "Internet service classes under competition," *IEEE J. Sel. Areas Commun.*, vol. 18, no. 2, pp. 2490C2498, Feb. 2000.

[51] R. Jain, T. Mullen, and R. Hausman, "Analysis of Paris Metro pricing strategy for QoS with a single service provider," *IEEE/IFIP Intl. Workshop Quality of Service*, 2001.

[52] *AT&T Wireless Data Plan*, available at http://www.wireless.att.com.

[53] *HSPA to LTE-advanced: 3GPP broadband evolution to IMT-advanced*, Rysavy Research for 3G Americas, Sep. 2009.

[54] *Mobile users wary of usage-based pricing: Survey*, Reuters, Nov. 18, 2010. http://www.reuters.com/article/idUSN1827905920101118.

[55] R. Johari, G.Y. Weintraub, and B. Van Roy, "Investment and market structure in industries with congestion," *Operations Research*, vol. 58 no. 5, pp. 1303-1317, 2010.

[56] G. W. Evans and S. Honkapohja, *Learning and Expectations in Macroeconomics*, Princeton, NJ: Princeton Univ. Press, 2001.

[57] D. P. Bertsekas and J. N. Tsitsiklis, *Parallel and Distributed Computation: Numerical Methods*, Belmont, MA: Athena Scientific, 1989.

[58] A. Gosh and P. McAfee, "Incentivizing high-quality user-generated content," *20th Intl. Conf. World Wide Web*, 2011.

[59] S. Jain, Y. Chen, and D. C. Parkes, "Designing incentives for online question and answer forums," *ACM Conf. Electronic Commerce*, 2009.

[60] V. K. Singh, R. Jain, and M. S. Kankanhalli, "Motivating contributors in social media networks," *ACM SIGMM Workshop on Social Media*, 2009.

[61] J. Park and M. van der Schaar, "A game theoretic analysis of incentives in content production and sharing over peer-to-peer networks," *IEEE J. Sel. Signal Process.*, vol. 4, no. 4, pp. 704-717, Aug. 2010.

[62] S. Ren and M. van der Schaar, "Pricing and distributed power control in wireless relay networks," *IEEE Trans. Signal Process.*, vol. 59, no. 6, pp. 2913-2926, June 2011.

[63] A. Hagiu, "Merchant or two-sided platform?" *Review of Network Economics*, vol. 6, no. 2, pp. 115-133, Jun. 2007.

[64] J. Musacchio, G. Schwartz, and J. Walrand, "A two-sided market analysis of provider investment incentives with an application to the net-neutrality issue," *Review of Network Economics*, vol. 8, no. 1, pp. 22-39, Mar. 2009.

[65] O. Candogan, K. Bimpikis, and A. Ozdaglar, "Optimal pricing in networks with externalities," *submitted* for publication, 2010.

[66] Y. Jin, S. Sen, R. Guerin, K. Hosanagar, and Z.-L. Zhang, "Dynamics of competition between incumbent and emerging network technologies," *NetEcon*, Aug. 2008.

[67] S. Shakkottai, R. Srikant, A. Ozdaglar, and D. Acemoglu, "The price of simplicity," *IEEE J. Sel. Areas in Commun.*, vol. 26, no. 7, pp. 1269-1276, Sep.2008.

[68] A. K. Dixit and J. E. Stiglitz, "Monopolistic competition and optimum product diversity," *American Economic Review*, vol. 67, no. 3, pp. 297-308, 1977.

[69] J. R. Munkres, *Elements of Algebraic Topology*, New York: Perseus Books Pub., 1993.

[70] http://en.wikipedia.org/wiki/User-generated_content

[71] M. Armbrust, A. Fox, R. Griffith, A. D. Joseph, R. H. Katz, A. Konwinski, G. Lee, D. A. Patterson, A. Rabkin, I. Stoica, and M. Zaharia, "Above the clouds: A Berkeley view of cloud computing," *UC Berkeley Tech. Report UCB/EECS-2009-28*, Feb. 2009.

[72] AT&T Cloud Services, http://www.synaptic.att.com

[73] R. Ferzli and I. Khalife, "Mobile cloud computing educational tool for image/video processing algorithms," *IEEE DSP/SPE*, 2011.

[74] X. Geng and A. B. Whinston, "Profiting from Value-Added Wireless Services," *Computer*, vol. 34, no. 8, pp. 87-89, Aug. 2001.

[75] D. Niyato and E. Hossain, "Competitive pricing for spectrum sharing in cognitive radio networks: Dynamic game, inefficiency of Nash equilibrium, and collusion," vol. 26, no. 1, pp. 192-202, Jan. 2008.

[76] Y. Xi and E. Yeh, "Pricing, competition, and routing in relay networks," *Allerton*, 2009.

[77] S. Ren and M. van der Schaar, "Pricing and distributed power control in wireless relay networks," *IEEE Trans. Signal Process.*, vol. 59, no. 6, pp. 2913-2926, Jun. 2011.

[78] S. Kosta, A. Aucinas, P. Hui, R. Mortier, and X. Zhang, "ThinkAir: Dynamic resource allocation and parallel execution in the cloud for mobile code offloading," *IEEE Infocom*, 2012.

[79] Y. Wen, W. Zhang, and H. Luo, "Energy-pptimal mobile application execution: Taming resource-poor mobile devices with cloud clones," *IEEE Infocom*, 2012.

[80] K, Kant, "Data center evolution: A tutorial on state of the art, issues, and challenges," *Elsevier Networks*, vol. 53, no. 17, Dec. 2009.

[81] J. Choi, S. Govindan, B. Urgaonkar, and A. Sivasubramaniam, "Power consumption prediction and power-aware packing in consolidated environments," *IEEE Trans. Computers*, vol. 59, no. 12, pp. 1640-1654, Dec. 2010.

[82] X. Fan, W.-D. Weber, and L. A. Barroso, "Power provisioning for a warehouse-sized computer," *ACM ISCA*, 2007.

[83] A. Gandhi, M. Harchol-Balter, and C. L. R. Das, "Optimal power allocation in server farms," *ACM Sigmetrics,* 2009.

[84] C. E. Bash, C. D. Patel, and R. K. Sharma, "Dynamic thermal management of aircooled data centers," *ITHERM*, 2006.

[85] Z. Wang, A. McReynolds, C. Felix, C. Bash, and C. Hoover, "Kratos: Automated management of cooling capacity in data centers with adaptive vent tiles," *IMECE*, 2009.

[86] C. Patel, R. Sharma, C. Bash, and A. Beitelmal, "Energy flow in the information technology stack," *IMECE*, 2006.

[87] M. A. Adnan, Y. Ma, R. Sugihara, and R. Gupta, "Dynamic deferral of workload for capacity provisioning in data centers," `http://arxiv.org/abs/1109.3839`.

[88] M. Lin, A. Wierman, L. L. H. Andrew, and E. Thereska, "Dynamic right-sizing for power-proportional data centers," *IEEE Infocom*, 2011.

[89] B. Guenter, N. Jain, and C. Williams, "Managing cost, performance and reliability tradeoffs for energy-aware server provisioning," *IEEE Infocom*, 2011.

[90] A. Qureshi, R. Weber, H. Balakrishnan, J. Guttag, and B. Maggs, "Cutting the electric bill for internet-scale systems," *ACM Sigcomm*, 2009.

[91] Z. Liu, M. Lin, A. Wierman, S. Low, and L. H. Andrew, "Greening geographical load balancing", *ACM Sigmetrics*, 2011.

[92] L. Rao, X. Liu, L. Xie, and Wenyu Liu, "Reducing electricity cost: optimization of distributed Internet data centers in a multi-electricity-market environment," *IEEE Infocom*, 2010.

[93] D. Xu and X. Liu, "Geographic trough filling for Internet datacenters," *IEEE Infocom*, 2012.

[94] Y. Yao, L. Huang, A. Sharma, L. Golubchik, and M. J. Neely, "Data centers power reduction: A two time scale approach for delay tolerant workloads," *IEEE Infocom*, 2012.

[95] Y. Guo, Z. Ding, Y. Fang, and D. Wu, "Cutting down electricity cost in Internet data centers by using energy storage," *IEEE Globecom*, 2011.

[96] Z. Liu, Y. Chen, C. Bash, A. Wierman, D. Gmach, Z. Wang, M. Marwah, and C. Hyser, "Renewable and cooling aware workload management for sustainable data centers," *ACM Sigmetrics*, 2012.

[97] *Electricity Energy Storage Technology Options: A White Paper Primer on Applications, Costs and Benefits,* 2010, (`http://energy.gov/`).

[98] I. C. Paschalidis and J. N. Tsitsiklis, "Congestion-dependent pricing of network services," *IEEE/ACM Trans. Networking*, vol. 8, no. 2, pp. 171-184, Apr. 2000.

[99] P. Hande, M. Chiang, A. R. Calderbank, and S. Rangan, "Network pricing and rate allocation with content provider participation," *IEEE Infocom*, Apr. 2009.

[100] H. Mohsenian-Rad, V. W.S. Wong, J. Jatskevich, R. Schober, and A. Leon-Garcia, "Autonomous demand side management based on game-theoretic energy consumption scheduling for the future smart grid," *IEEE Trans. Smart Grid*, vol. 1, no. 3, pp. 320-331, Dec. 2010.

[101] I. Menache, A. Ozdaglar, and N. Shimkin, "Socially optimal pricing of cloud computing resources," *ValueTools*, 2011.

[102] Q. Wang, K. Ren, and X. Meng, "When cloud meets eBay: Towards effective pricing for cloud computing," *IEEE Infocom*, 2012.

[103] R. Pal and P. Hui, "Economic models for cloud service markets pricing and capacity planning," *Elsevier Theoretical Computer Science*, accepte and to appear.

[104] C. Maglaras and A. Zeevi, "Pricing and capacity sizing for systems with shared resources: Approximate solutions and scaling Relations," *Management Science*, vol. 49, no. 8, pp. 1018-1038, Aug. 2003.

[105] K. Charnsirisakskul, P. Griffin, and P. Keskinocak, "Pricing and scheduling decisions with leadtime flexibility," *Eur. J. Oper. Res.*, vol. 171, no. 1, pp. 153C169, 2006.

[106] Z.-L. Chen and N. G. Hall, "The coordination of pricing and scheduling decisions," *Manufacturing and Service Operations Management*. vol. 12, no. 1, pp. 77-92, Winter 2010.

[107] J. C. Hallak, "The effect of cross-country differences in product quality on the direction of international trade 2002," *Working Paper*, Univ. Michigan, Ann Arbor, MI. (http://ideas.repec.org/p/mie/wpaper/493.html)

[108] J. Nair, A. Wierman, and B. Zwart, "Exploiting network effects in the provisioning of large scale systems," *IFIP Performance*, 2011.

[109] M.D. Ilic, L. Xie and J. Joo, "Efficient coordination of wind power and price-responsive demand Part I: theoretical foundations," *IEEE Trans. Power Systems*, vol. 26, no. 4, pp. 1875-1884, Nov. 2011.

[110] M. J. Neely, "Universal scheduling for networks with arbitrary traffic, channels, and mobility," *Technical Report*, 2010.

[111] J. C. Rochet and J. Tirole, "Two-sided markets: A progress report ," *RAND J. Economics*, vol. 37, no. 3, pp. 645-667, Autumn 2006.

[112] `http://www.caiso.com`

[113] S. Boyd and L. Vandenberghe, *Convex Optimization*, Cambridge University Press, 2004.