

Profit-Driven Dynamic Cloud Pricing for Multiserver Systems Considering User Perceived Value

Peijin Cong, Liying Li, Junlong Zhou, Kun Cao, Tongquan Wei, Mingsong Chen, and Shiyan Hu

Abstract—With the rapid deployment of cloud computing infrastructures, understanding the economics of cloud computing has become a pressing issue for cloud service providers. However, existing pricing models rarely consider the dynamic interaction between user requests and the cloud service provider, thus can not accurately reflect the law of supply and demand in marketing. In this paper, we first propose a dynamic pricing model based on the concept of user perceived value in the domain of economics that accurately captures the real supply and demand situation in the cloud service market. Subsequently, we design a profit maximization scheme based on the dynamic pricing model that optimizes profit of the cloud service provider without violating user service-level agreement. Finally, we propose a dynamic closed loop control scheme to adapt the cloud service price and multiserver configurations to the changes in cloud computing environment. Extensive experiments using data extracted from real-world applications validate the effectiveness of the proposed user perceived value-based pricing model and the dynamic profit maximization scheme. Our proposed profit maximization scheme achieves 31.32% and 22.76% more profit compared to two state of the art benchmarking methods, respectively.

Index Terms—Cloud computing, dynamic pricing model, user perceived value, profit maximization, closed loop control.

1 INTRODUCTION

CLOUD computing has become an effective commercial computing model that distributes user requests on a pool of servers and delivers hosted services over Internet. As a business model, it turns resources of computing, storage, and communication into ordinary commodities and utilities in a pay-as-you-go manner [1], [2], [3], [4], [5]. It is natural for cloud service providers to pursue the goal of profit maximization. Thus, the cloud service pricing strategy is of particular importance to cloud service providers.

The pricing model of a cloud service provider in cloud computing consists of two parts, namely, the revenue and the cost [6]. From the perspective of a cloud service provider, the revenue is the income that the cloud service provider has from the sale of cloud services to users. The cost is the expenditure that includes not only the rental and electricity fees to operate multiserver systems, but also the reward and penalty paid by the cloud service provider to users based on service-level agreement. Profit maximization can be achieved by increasing revenue or reducing cost. On one hand, cloud service providers attempt to increase revenue by setting a high price for cloud services and attracting a great amount of service purchases. However, service price and purchase amount interplay and cannot be optimized simultaneously [7]. On the other hand, cloud service providers try to reduce operating cost. Thus, aspects such as multiserver configurations and electricity price should be considered in cloud pricing modeling.

Numerous investigations have been made into pricing mechanisms for profit maximization in cloud computing. Fixed pricing strategies such as pay-per-use, subscription based pricing, and tiered pricing are the most common pricing methods used by major cloud service providers [8], [9], [10]. For example, Li [8] proposed a flat rate pricing strategy that sets a fixed price for all service requests. Kesidis et al.

[9] pointed out that usage-based pricing strategy can use cloud resources more efficiently when compared with flat rate pricing strategy. However, these fixed pricing methods cannot meet the dynamic needs of users and cannot reflect the market situation of supply and demand.

The disadvantages of fixed pricing strategies necessitate dynamic pricing strategies that adjust price of cloud services according to market situations and user requirements for service quality. Macias et al. [11] proposed a genetic model based dynamic pricing strategy that obtains optimal pricing in an iterative way. This strategy offers competitive prices in the negotiation of services in cloud computing markets. Amazon [12], [13] utilized a spot pricing strategy that dynamically adjusts prices for a virtual service instance to accommodate changes in supply and demand. Based on a study of the spot price history of Amazon, Xu et al. [14] proposed a dynamic pricing strategy to better understand the current market demand. Zhao et al. [15] designed an efficient online algorithm for dynamic pricing of virtual machine resources across datacenters in a geo-distributed cloud to pursue long-term profit maximization. Although these works investigate dynamic pricing strategies from different perspectives, the service-level agreement is not considered in the presented pricing mechanisms.

A service-level agreement is defined as an official commitment that prevails between a service provider and a client [16]. It uses a price compensation mechanism that gives users certain compensations when their service requests are processed with low quality of service. Cao et al. [6] presented a pricing model that takes the service-level agreement and the consumer satisfaction into considerations to maximize the profit of cloud service providers. Ghamkhari et al. [17] proposed a two-tier ladder-like charging method to ensure user satisfaction. A cloud service

provider will charge users if user requests are processed before deadlines, the cloud service is free otherwise. Lee and Irwin [18], [19] et al. claimed that the price of the cloud service will decrease as the waiting time of service requests grows until the cloud service is free. These works study the service-level agreement to ensure user satisfaction in the pricing process for profit maximization. However, they ignore the crucial concept of user perceived value in traditional market environment, which reflects the users' willingness to purchase cloud services and ultimately has influence on the profit of cloud service providers.

In this paper, we propose a user perceived value-based dynamic pricing mechanism that conforms to the law of supply and demand in economics. The novel contributions of this paper are summarized as follows:

- We propose a dynamic pricing model that considers the interaction between cloud users and the cloud service provider. The model built upon the concept of user perceived value, user reward, and cloud service provider penalty in the domain of economics accurately captures the dynamics of supply and demand in cloud pricing strategies. In particular, user perceived value is nicely modeled using kernel density estimation method.
- We propose a profit maximization scheme based on the dynamic pricing model to optimize the profit of the cloud service provider by configuring multi-server systems under the constraint of service-level agreement. The proposed scheme includes a runtime control loop that specifically considers the dynamic cloud computing environment such as fluctuating electricity bill and rental fees that are not integrated in the pricing model.
- Extensive simulation experiments show that the proposed scheme not only follows the supply and demand law in market, but also is superior to the state of the art benchmarking cloud pricing mechanisms. The proposed scheme achieves 31.32% and 22.76% more profit as compared to two state of the art benchmarking methods, respectively.

The remainder of the paper is organized as follows. Section 2 presents the system architecture and models, Section 3 presents the problem definition and overview of the proposed scheme. Section 4 describes the proposed user perceived value-based pricing mechanism. The effectiveness of the proposed scheme is validated in Section 5 and concluding remarks are given in Section 6.

2 SYSTEM ARCHITECTURE AND MODELS

We consider a common three-tier cloud service provision structure that consists of cloud users, cloud service providers, and cloud infrastructure vendors [6], [18], [20]. Among the three entities that form a market in cloud computing, the infrastructure vendor charges the cloud service provider for renting infrastructures to deploy service capacity, and the cloud service provider charges cloud users for processing their service requests. In this paper, cloud users and the cloud service provider are of our particular interest. We introduce our cloud user model and cloud service provider model in the following subsections.

2.1 Cloud User Model

To maximize the profit of a cloud service provider, the cloud service provider needs to know the aggregate demands of all users. When a cloud service provider sets up the price of a service, different users have different responses to this price. In this paper, we propose a user perceived value oriented pricing strategy for cloud computing services. In this subsection, we introduce the concepts of user perceived value and user request (or demand) distribution.

2.1.1 User Perceived Value

In conventional markets, the arrival rate of customers to a store is often a response to their regular buying patterns rather than a reaction to individual prices [7]. Thus, it is reasonable to assume that the change of the list price has no effect on the total number of customers who are visiting the store. Typically, not all of the customers are willing to buy a specific commodity. That is, the total number of customers who buy commodities are no larger than the total number of people that visit the store.

Customer perceived value is the fundamental basis for all marketing activities [21]. It reflects the worth that a product or service has in the mind of a consumer and has been widely used in modeling other markets [22]. In general, customers are unaware of the true cost of production for the products they buy, instead, they simply have an internal feeling for how much certain products are worth to them. In the conventional market environment, only the customer whose perceived value is higher than the real price of the product is willing to pay for the product. As with conventional market commodities, the cloud computing service can be seen as a special commodity which follows market rules. That is, the price of cloud computing service is also dictated by the supply and demand in the market.

In this paper, we use X_i to denote the perceived value that user i has for the service. X_i is a continuous random variable and $0 \leq X_i < \infty$ holds. As with other benchmarking pricing models [23], X_1, X_2, \dots, X_n are assumed to be independent and identical random variables. The probability density function of the perceived value X , denoted by $f(x)$, is known or can be estimated a priori. Perceived value is a process of valuing and is much harder to determine. Roig et al. [24] observed that the customer value is perceived by customers, and cannot be determined objectively by the seller. Factors such as scarcity, marketing efforts, novelty, and brand associations all play into customer perceived value [25]. Usually, consumers will offer a range of price options. Thus, in the experimental section, we use a normal distribution to describe the initial distribution of user perceived value. Subsequently, we fit the probability distribution using kernel density estimation based on historical price data. Kernel density estimation is a non-parametric way to estimate the probability density function of a random variable based on a finite data sample [26].

In the following sections, we adopt the terminology of customer perceived value used in traditional market environment. We take the cloud computing environment as a store and the cloud computing service is deemed as a special commodity provided in the store. The terminology of customer perceived value and user perceived value are used interchangeably.

2.1.2 User Demand Distribution

Unlike traditional methods that use the expected demand to model user behavior [23], [27], we use the probability distribution of the total demands in this work to model user service requests. We consider a slotted time model that deals with the pricing decision and constraints for sales periods T of equal length. Let τ denote the length of each time slot over the sales period T , and N be the number of time slots τ over the sales period T . That is, $T = N \times \tau$. A cloud service provider sets list price for the service at the beginning of regular sales periods. The list price during each sales period is assumed to be constant, but varies from period to period.

Suppose that the cloud service provider will charge ω per user for a specific cloud service during a sales period T . Let n denote the total number of users that have interest in the service at the price of ω during the sales period T , and λ_u denote the number of users arriving per unit time, respectively. n is assumed to be independent of all other parameters of the system, and is a discrete Poisson random variable distributed as [27]

$$P(n|\lambda_u) = \frac{(\lambda_u T)^n e^{-\lambda_u T}}{n!}, \quad n = 0, 1, 2, \dots, \infty. \quad (1)$$

The user arrival rate λ_u may not be constant in many situations. Taking into account the heterogeneity of arrival rate, a Gamma distribution characterized by parameters (α, β) is utilized to represent the arrival rate λ_u , the probability density function of which is given by

$$g(\lambda_u) = \frac{1}{\Gamma(\alpha)\beta^\alpha} \lambda_u^{(\alpha-1)} e^{-\lambda_u/\beta}, \quad 0 \leq \lambda_u \leq \infty, \quad (2)$$

where the expectation and variance of λ_u is given by $E[\lambda_u] = \alpha\beta$ and $Var[\lambda_u] = \alpha\beta^2$, respectively, and $\Gamma(\alpha)$ is a complete gamma function.

Among the n users, any one whose perceived value of the service is no less than the list price ω is considered as a potential buyer of the service. Let m denote the number of potential buyers. It is a non-negative discrete random variable taking the value of $0, 1, 2, \dots, \infty$ and $m \leq n$ holds. Based on user perceived value, in this subsection, we use $F(\omega)$ to represent the cumulative distribution function of x evaluated at ω . The $F(\omega)$ is a non-decreasing function of ω , and $0 \leq F(\omega) \leq 1$ and $\lim_{\omega \rightarrow \infty} F(\omega) = 1$ hold [28]. Let $P_\omega(m|n)$ indicate the probability that m out of n users are inclined to buy in the sales period when the service price is set equal to ω . It follows a binomial distribution of probability, which is given by

$$P_\omega(m|n) = \binom{n}{m} [1 - F(\omega)]^m [F(\omega)]^{(n-m)}. \quad (3)$$

Combining (1)-(3), we can derive the probability of having m potential buyers during the sales period T when the service price is set equal to ω . The probability is denoted by $P_\omega(m)$ and given by

$$\begin{aligned} P_\omega(m) &= \int_{\lambda_u=0}^{\infty} \sum_{n=0}^{\infty} P_\omega(m|n) P(n|\lambda_u) g(\lambda_u) d\lambda_u \\ &= \binom{m+\alpha-1}{m} \left[\frac{\beta T [1 - F(\omega)]}{1 + \beta T [1 - F(\omega)]} \right]^m \left[\frac{1}{1 + \beta T [1 - F(\omega)]} \right]^\alpha. \end{aligned} \quad (4)$$

Clearly, it is a negative binomial distribution. As a result, the expected number of actual buyers of the service at price ω during sales period T , which is denoted by $E_\omega(m)$, can be calculated as

$$E_\omega(m) = \alpha\beta T (1 - F(\omega)), \quad (5)$$

where α and β are parameters of the Gamma distribution of user arrival rate λ_u , and $F(\omega)$ is the cumulative distribution function of x evaluated at ω . The revenue of the cloud service provider in a sales period T is thus given by

$$Revenue = \omega \times E_\omega(m) = \omega\alpha\beta T (1 - F(\omega)). \quad (6)$$

2.2 Cloud Service Provider Model

The cloud service provider rents a multiserver system that is constructed and maintained by an infrastructure vendor to serve user service requests. The architecture details of the multiserver system are quite flexible [6]. They can be blade centers where each server is a server blade [29], clusters of traditional servers where each server is an ordinary processor [30], and multicore server processors where each server is a single core [31]. For the ease of presentation, these blades/processors/cores are simply called servers. Users submit their service requests to the cloud service provider, and the cloud service provider serves these service requests (i.e., run these tasks) on the multiserver system.

2.2.1 Multi-Server Model

We consider a multiserver system that consists of M homogeneous servers operating at a common speed of s . The multiserver system can be modeled as an M/M/M queuing system where arrivals of user service requests governed by a Poisson process form a single queue and M servers can process these service requests in parallel. Let ρ be the service rate of user service requests that arrive at the rate of λ_u . It is clear that ρ user service requests can be processed by servers if the number of user service requests in the system is not greater than M . The service time of a user service request on a server is an exponential random variable denoted by $x_1 = r/s$ with mean $\bar{x}_1 = \bar{r}/s$, where r is the number of instructions to be executed for the service request. A first-come-first-served (FCFS) queue of infinite capacity is maintained by the multiserver system for waiting tasks when all the servers are busy. Let ρ be server utilization, which is defined as the average percentage of time that a server is busy. It can be expressed as

$$\rho = \frac{\lambda_u}{M\rho} = \frac{\lambda_u}{M\frac{\bar{r}}{s}} = \frac{\lambda_u \bar{r}}{Ms}. \quad (7)$$

Let P_k be the probability of k service requests being waiting or processing in the M/M/M queuing system. Based on queuing theory [6], [32], [33], P_k is given by

$$P_k = \begin{cases} P_0 \frac{(M\rho)^k}{k!}, & k \leq M \\ P_0 \frac{M^M \rho^k}{M!}, & k \geq M \end{cases}, \quad (8)$$

where P_0 is the probability that there are no tasks in the queue, and is formulated into [32]

$$P_0 = \left(\sum_{k=0}^{M-1} \frac{(M\rho)^k}{k!} + \frac{(M\rho)^M}{M!} \cdot \frac{1}{1-\rho} \right)^{-1}.$$

The probability that there are exact M service requests in the system is thus given by $P_M = P_0 \frac{(M\rho)^M}{M!}$. Through Taylor series expansions of $\sum_{k=0}^{M-1} \frac{(M\rho)^k}{k!} \approx e^{M\rho}$ and $M! \approx \sqrt{2\pi M} \left(\frac{M}{e}\right)^M$, it can be rewritten as

$$P_M = \frac{1-\rho}{\sqrt{2\pi M}(1-\rho)\left(\frac{e^{\rho-1}}{\rho}\right)^{M+1}}. \quad (9)$$

This form of P_M is necessary for deriving multiserver configurations in Section 4.

When all the servers in the system are busy, a newly submitted service request must wait and will be inserted into the FCFS queue. Let P_q denote the probability of queuing a newly arrived task when no servers are idle at the time of arrival. P_q can be formulated as

$$P_q = \sum_{k=M}^{\infty} P_k = \frac{P_M}{1-\rho}. \quad (10)$$

Let \bar{N} be the average number of requests being waiting or executing in the multiserver system. \bar{N} is calculated as

$$\bar{N} = \sum_{k=0}^{\infty} kP_k = M\rho + \frac{\rho}{1-\rho}P_q. \quad (11)$$

The average service response time \bar{R} is defined as the time elapsed between the time when a service request is submitted and the time when the service request is finished. In this paper, it is adopted to evaluate the service quality. It is in fact the sum of task execution time and waiting time, and can be derived by applying Little's Law [34] as

$$\bar{R} = \frac{\bar{N}}{\lambda_u} = \frac{\bar{N}}{\lambda_u} \left(1 + \frac{P_q}{M(1-\rho)}\right) = \frac{\bar{N}}{\lambda_u} \left(1 + \frac{P_M}{M(1-\rho)^2}\right). \quad (12)$$

The average service response time \bar{R} is utilized in this paper as a metric for service-level agreement. If the response time of a service exceeds the predefined deadline, the service-level agreement is deemed to be violated.

2.2.2 Bill and Rent

A cloud service provider needs to rent infrastructure and pay electricity to maintain the operation of the computing infrastructure. Let δ be the fee the cloud service provider pays to rent a server per second during a sales period T , the rent the cloud service provider needs to pay for a system of M servers during the sales period T is

$$Rent = M\delta \times T. \quad (13)$$

As a portion of the cloud service cost, electricity fee has become a significant expense for today's data centers. It can be derived by multiplying energy consumed by a server with electricity price. The energy consumed by a

server can be modeled at different levels of abstraction. At the abstraction level of digital CMOS circuit, the power consumption, which is denoted by P_{tot} , can be modeled as

$$P_{tot} = P_{sta} + P_{dyn}, \quad (14)$$

where P_{sta} is the static power dissipation while P_{dyn} is the dynamic power dissipation. P_{sta} is independent of switching activity and maintains the basic circuit state, thus can be deemed as a constant [6]. P_{dyn} is related to processor switching activity and dominates the total power consumption, which can be formulated as a function of supply voltage v and processing speed s . In addition, the supply voltage is usually linearly proportional to the processing speed, i.e., $v \propto s$. The dynamic power consumption P_{dyn} is then expressed as ξs^γ , where ξ is a processor dependent coefficient and γ is a constant that equals to $2\phi + 1$ ($\phi > 0$). Based on the static and dynamic power consumption described above, we use the following Equation (15) to denote the total power consumption of a multiserver system, which is,

$$P_{tot} = M((P_{dyn} - P_{sta})\rho + P_{sta}), \quad (15)$$

where M is the number of servers and ρ is the server utilization.

Let E^T denote the energy consumed by all M servers in the system during the sales period T . It is given by

$$E^T = M((P_{dyn} - P_{sta})\rho + P_{sta}) \times T. \quad (16)$$

Let $C^T(E^T)$ denote the price of the energy consumed by all servers in the sales period T , then $C^T(E^T)$ can be formulated as

$$C^T(E^T) = \begin{cases} k_1^T, & 0 \leq E^T \leq l_{th}^T \\ k_2^T, & E^T > l_{th}^T \end{cases} \quad (17)$$

where $k_1^T, k_2^T > 0$ are differentiated price and l_{th}^T is the energy consumption threshold in the sales period T . The electricity bill of the multiserver system in the sales period T is hence formulated as

$$\begin{aligned} Bill &= E^T \times C^T(E^T) \\ &= M((P_{dyn} - P_{sta})\rho + P_{sta}) \times T \times C^T(E^T). \end{aligned} \quad (18)$$

2.3 Reward and Penalty Mechanism

Oftentimes, users have different sensitivities to postponing their requests. For users whose service requests can be deferred to a certain extent, the cloud service provider will reward them based on the degree of deferment. However, once the deferment of service requests exceeds a threshold, the cloud service provider will compensate users based on the degree of the deferment. In the following, we will discuss the reward and penalty mechanism from perspectives of users and the cloud service provider, respectively.

2.3.1 User Reward Model

We divide users into I types, each type of users has a sensitivity to the service deferment of their service requests. We define a sensitivity factor, denoted by ψ_i , to quantify the

sensitivity of users of type i to the deferment of their service requests, which is denoted by D_i . For the users of type i , the factor ψ_i is negatively proportional to the sensitivity of the users to the deferment of their service requests. That is to say, a larger sensitivity factor indicates a more delay-sensitive user service request. For users running interactive applications with no delay tolerance, we set $\psi_i = \infty$.

The cloud service provider will return more rewards to those users who are less sensitive to service deferment. Let L_i be the monetary loss of the users of type i due to their degree of sensitivity to service deferment. The users with a larger degree of sensitivity (ψ_i) to service deferment (D_i) will get less rewards from the cloud service provider, resulting in a larger amount of monetary loss (L_i) of the users of type i . The monetary loss function of the users of type i is given by $L_i = \psi_i D_i$.

We define a reward function, denoted by h_i , to represent the reward the cloud service provider returns to users of type i . The reward h_i is a function of the service deferment D_i , and is given by $h_i = \theta \cdot \log(1 + D_i)$, where $\theta > 0$ and $0 \leq D_i \leq D_{max}$ hold. θ is called the reward factor and D_{max} is the maximum value of service deferment.

Users need to make decisions on their own service deferment D_i to get the maximum monetary reward from the cloud service provider based on the monetary loss and reward function. Thus, the optimization problem is to maximize $(h_i - L_i)$ subject to $(0 \leq D_i \leq D_{max})$, where D_i is regarded as a continuous variable to simplify the optimization problem, and the solution to the problem is given by

$$D_i = \max(\min(\frac{\theta}{\psi_i} - 1, D_{max}), 0). \quad (19)$$

Next, substitute D_i back into the reward function h_i , we have

$$h_i = \theta \log(1 + \max(\min(\frac{\theta}{\psi_i} - 1, D_{max}), 0)). \quad (20)$$

Let *Reward* denote the total monetary reward returned to users from the cloud service provider over the sales period T , then we have

$$Reward = \sum_{N'=1}^N \sum_{i=1}^I (h_i - L_i) \lambda_u^i [N'] \tau, \quad (21)$$

where $\lambda_u^i [N']$ denotes the arrival rate of the users of type i in the N' th time slot. In practice, the cloud service provider can learn the sensitivity factor ψ_i from experiments or historical data. The adopted user reward model is similar to the one presented in [35].

2.3.2 Cloud Service Provider Penalty Model

The high degree of user satisfaction is determined by the fast response of a cloud service provider to users' service requests. Once the service response time exceeds the threshold value specified in service-level agreement, users will be compensated by the cloud service provider for low quality of service. Let s_0 indicate the benchmarking speed of servers. Given server benchmarking speed (s_0), the average

response time of service requests (\bar{R}), and the number of instructions for each service request (r), the penalty function of the users of type i , denoted by u_i , can be formulated as

$$u_i(r, \bar{R}) = \begin{cases} 0, & 0 \leq \bar{R} \leq \frac{r}{s_0} + D_i \\ d(\bar{R} - \frac{r}{s_0} - D_i), & \frac{r}{s_0} + D_i < \bar{R} \leq (1 + \frac{\omega}{d}) \frac{r}{s_0} + D_i \\ \omega, & \bar{R} > (1 + \frac{\omega}{d}) \frac{r}{s_0} + D_i \end{cases} \quad (22)$$

where d is the degree of punishment suffered when the service-level agreement is violated. D_i denotes the service deferment of the requests of the users of type i , and ω is the price of the cloud service charged by the cloud service provider to the users.

The details of Equation (22) are described below. For the users of type i , if the average response time satisfies $0 \leq \bar{R} \leq \frac{r}{s_0} + D_i$, the cloud service provider will regard this execution of the service request as a successful process with high quality of service and the users will not be compensated by the cloud service provider. Otherwise, if the average response time satisfies $\frac{r}{s_0} + D_i < \bar{R} \leq (1 + \frac{\omega}{d}) \frac{r}{s_0} + D_i$, the cloud service provider will regard this execution of the service request as a process with low quality of service. In this case, the compensation provided by the cloud service provider to the users will increase linearly as the average response time \bar{R} increases. Finally, if the average response time satisfies $\bar{R} > (1 + \frac{\omega}{d}) \frac{r}{s_0} + D_i$, the cloud service provider will regard this execution of the service request as a failed process and will not charge the users for this execution.

We use *Penalty* to denote the total compensation provided by the cloud service provider to users, then we have

$$Penalty = \sum_{N'=1}^N \sum_{i=1}^I u_i(r, \bar{R}) \lambda_u^i [N'] \tau, \quad (23)$$

where $\lambda_u^i [N'] \tau$ is the average number of user requests of type i in the time slot τ , and $u_i(r, \bar{R})$ is the incurred penalty for the service requests due to low quality of service.

2.3.3 Gross Profit

The gross profit a cloud service provider earns is the total revenue subtracted by the cost of generating that revenue. In other words, gross profit is sales minus cost of the cloud service sold. Assuming the price of cloud service is constant in a sales period T , the revenue earned is given by $\omega \cdot E_\omega(m)$, where ω denotes the service price per user and $E_\omega(m)$ indicates the expected number of actual buyers at price ω during the sales period T .

Besides the reward for flexible users and penalty for low quality of service mentioned above, the cost of cloud service provider sold also consists of the cost paid to rent cloud computing infrastructure, and the electricity expense incurred by the cloud service provider to maintain the operation of the computing infrastructure. We define the profit of the cloud service provider in a sales period T as the revenue minus the various expenses including the reward cost, penalty cost, electricity cost, and rental cost, that is,

$$Profit = Revenue - Reward - Penalty - Bill - Rent, \quad (24)$$

where *Revenue*, *Reward*, *Penalty*, *Bill*, and *Rent* are given in Equations (6), (21), (23), (18), and (13), respectively.

3 PROBLEM DEFINITION AND OVERVIEW OF THE PROPOSED APPROACH

In this section, we formally define the profit maximization problem, followed by a brief overview of our proposed approach to the profit maximization problem.

3.1 Problem Definition

The price of a cloud service interplays with users who purchase the service, which in turn affects the revenue of the cloud service provider. This paper aims to maximize the profit of the cloud service provider by deriving the optimal number of servers, operating speed of servers, and price of cloud services provided without violating the user service-level agreement. We assume that the cloud service provider optimizes its decisions at the beginning of each sales period T . Let b_1 denote the upper bound on the power consumption of the M servers, and b_2 be the upper bound on the expected response time of user requests. The optimization problem we will solve is thus formulated into

$$\text{Maximize: } Profit \quad (25)$$

$$\text{subject to: } \theta \geq 0 \quad (26)$$

$$P_{tot} \leq b_1 \quad (27)$$

$$\bar{R} \leq b_2 \quad (28)$$

$$0 \leq \phi_i[N'] \leq \lambda_u^i[N']\tau, \quad \forall i \in I, N' \in N \quad (29)$$

$$\sum_{N'=1}^{N+|D_i|} (\lambda_u^i[N']\tau - \phi_i[N']) \geq \sum_{N'=1}^N \lambda_u^i[N']\tau, \quad \forall i \in I, N' \in N \quad (30)$$

where *Profit*, P_{tot} , and \bar{R} are given in Equations (24), (15), and (12), respectively.

In the above formulation, the reward factor θ is non-negative (Equation 26), the total energy consumption P_{tot} of the multiserver within the sales period T can not exceed b_1 (Equation 27), and the average service response time \bar{R} can not exceed b_2 (Equation 28). Equation (29) ensures that the amount of delayed service requests is nonnegative and not larger than the total number of service requests in each time slot τ , where $\phi_i[N']$ and $\lambda_u^i[N']\tau$ are the number of delayed and total service requests in the N' th time slot, respectively. Equation (30) ensures that the processing of the arrived user requests of type i in a sales period T can not be delayed longer than the allowed service deferment D_i of the user service requests. We will then use the augmented Lagrange multiplier method to solve the optimization problem, which will be described in detail in Section 4.

3.2 Overview of the Proposed Approach

The optimization problem given in Equation (25) tries to maximize the *Profit* of the cloud service provider under the constraints mentioned above. Figure 1 outlines the overview of our proposed approach to solve the optimization problem. We first establish user demand distribution based on the concept of user perceived value. Subsequently,

based on user demand distribution, we build revenue model and expenditure model to construct the profit maximization problem. Then, we propose to solve the optimization problem by using the augmented Lagrange multiplier method. Finally, since the parameters of electricity bill and rental fees will change over time, thus, the cloud computing environment is dynamic. We use a system monitor to check whether these environment parameters change, and propose a dynamic closed loop control scheme to adapt the service price and multiserver configurations to the changes in these parameters. The details of the proposed scheme are described in Section 4.

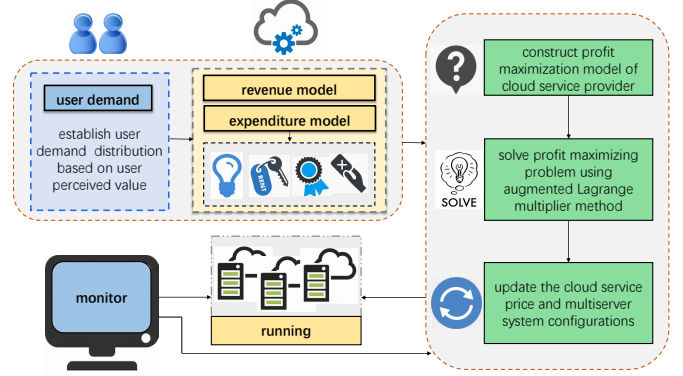


Figure 1: Overview of the proposed approach.

4 USER PERCEIVED VALUE-AWARE PROFIT OPTIMIZATION SCHEME

In this section, we first use augmented Lagrange multiplier method to solve the optimization problem and derive the optimal solution, including the optimal cloud service price, the number of servers, and the speed of servers. Subsequently, we propose a dynamic closed loop control scheme to adapt the service price and multiserver configurations to the dynamic cloud computing environment.

4.1 Create Augmented Lagrange Function

Unconstrained optimization can be solved in many ways, such as steepest descent method [36], Newton's method [37], multiplier method [38], etc. However, it is difficult to optimize constrained optimization directly. A common method to solve constrained optimization is to transform the constrained optimization problem into an unconstrained optimization problem. Numerous techniques on constrained optimization have been investigated in the literature [39], [40], [41]. Of these techniques, the augmented Lagrange multiplier method is a powerful tool for solving this class of problems, thus, is adopted in this work to solve the profit optimization problem given in Equation (25). We first create the augmented Lagrange function to transform the constrained optimization problem into an unconstrained optimization problem, and then use the multiplier method to solve the unconstrained optimization problem.

The *Bill* given in Equation (18) is a function of power consumption of the multiserver system, the length of the sales period T , and real-time price of electricity. Since real-time price is flat within each sales period T and T itself

is constant, the *Bill* for T is fixed and can be expressed as $Bill = b_3 P_{tot}$, where P_{tot} given in Equation (15) is the total power consumed by the multiserver system and b_3 is a constant coefficient. The profit optimization problem given in Equation (25) can then be re-written as

$$\left\{ \begin{array}{l} O(\omega, M, s) = \omega E_\omega[m] - b_3 P_{tot} - \delta MT \\ \quad - \sum_{N'=1}^N \sum_{i=1}^I (h_i - L_i) \lambda_u^i [N'] \tau \\ \quad - \sum_{N'=1}^N \sum_{i=1}^I u_i(r, \bar{R}) \lambda_u^i [N'] \tau \\ g_1(M, s) = b_2 - \bar{x}_1 \left(1 + \frac{P_M}{M(1-\rho)^2}\right) \geq 0 \\ g_2(M, s) = b_1 - M((\xi s^\gamma - P_{sta})\rho + P_{sta}) \geq 0 \end{array} \right. \quad (31)$$

where $O(\omega, M, s)$ denotes the objective function of *Profit* given in Equation (24), and $g_1(M, s)$ and $g_2(M, s)$ are constraint equations of M and s , respectively.

Next, we transform the problem given in Equation (31) with inequality constraints into an augmented Lagrange function. Let \mathbf{y} be the vector that converts the problem with inequality constraints to a problem with equality constraints, and \mathbf{v} be the Lagrange multiplier vector, the augmented Lagrange function is thus given by

$$\begin{aligned} \phi(\omega, M, s, \mathbf{y}, \mathbf{v}, \sigma) &= O(\omega, M, s) - \sum_{j=1}^2 v_j (g_j(M, s) - y_j^2) \\ &\quad + \frac{\sigma}{2} \sum_{j=1}^2 (g_j(M, s) - y_j^2)^2, \end{aligned} \quad (32)$$

where the constant parameter σ denotes the penalty factor and $\sigma > 0$ holds. The augmented Lagrange function given in Equation (32) can be converted into the form of

$$\begin{aligned} \phi(\omega, M, s, \mathbf{y}, \mathbf{v}, \sigma) &= O(\omega, M, s) \\ &\quad + \sum_{j=1}^2 \left[\frac{\sigma}{2} [y_j^2 - \frac{1}{\sigma} (\sigma g_j(M, s) - v_j)]^2 - \frac{v_j^2}{2\sigma} \right] \end{aligned} \quad (33)$$

by using the method of completing the square, a technique to derive the quadratic formula [42], and the function given in (33) can be easily maximized when

$$y_j^2 = \frac{1}{\sigma} \max(0, \sigma g_j(M, s) - v_j), j = 1, 2. \quad (34)$$

Plugging y_j^2 given in Equation (34) back into the original Formula (32), we have the desired augmented Lagrange function

$$\begin{aligned} \phi(\omega, M, s, \mathbf{v}, \sigma) &= O(\omega, M, s) \\ &\quad + \frac{1}{2\sigma} \sum_{j=1}^2 [(\max(0, v_j - \sigma g_j(M, s))]^2 - v_j^2]. \end{aligned} \quad (35)$$

Through this quadratic relaxation of the original problem given in Equation (31), we can derive analytical form of solutions to the profit maximization problem. We aim to maximize the profit of the cloud service provider and obtain the optimum solutions including service price ω , number of servers M , and speed of servers s . Specifically, we seek to solve the augmented Lagrange function given in Equation

(35) by first computing partial derivatives of Equation (31) with respect to ω , M and s . Here, we omit details (e.g., simple partial derivative process) of the derivation but only show key and difficult steps for solving the partial derivatives of Equation (31) with respect to ω , M and s below.

Calculate the partial derivative of Equation (31) with regard to ω : The partial derivative of $E_\omega[m]$ with regard to ω is $\frac{\partial E_\omega[m]}{\partial \omega} = -\alpha \beta t f(\omega)$, thus, the partial derivative of $O(\omega, M, s)$ with regard to ω is

$$\frac{\partial O(\omega, M, s)}{\partial \omega} = \frac{\partial E_\omega[m]}{\partial \omega} \cdot \omega + E_\omega[m] = \alpha \beta t (1 - \omega f(\omega) - F(\omega)),$$

where $f(\omega)$ and $F(\omega)$ are the probability density and the cumulative distribution function at ω , respectively.

Calculate the partial derivative of Equation (31) with regard to M : The partial derivative of $g_1(M, s)$ with regard to M can be expressed as

$$\begin{aligned} \frac{\partial g_1(M, s)}{\partial M} &= \frac{\partial}{\partial M} \left[-\frac{\bar{x}_1}{M} \left(\frac{1}{[\sqrt{2\pi M}(1-\rho)(e^\rho/e\rho)^M + 1](1-\rho)} \right) \right] \\ &\quad + \frac{\bar{x}_1}{M^2} \frac{P_M}{(1-\rho)^2}. \end{aligned} \quad (36)$$

Let $D_1 = \sqrt{2\pi M}(1-\rho)(e^\rho/e\rho)^M + 1 = \sqrt{2\pi M}(1-\rho)L + 1$, $D_2 = 1 - \rho$, and $L = (e^\rho/e\rho)^M$, then Equation (36) becomes

$$\frac{\partial g_1(M, s)}{\partial M} = \frac{\partial}{\partial M} \left[-\frac{\bar{x}_1}{M} \left(\frac{1}{D_1 D_2} \right) \right] + \frac{\bar{x}_1}{M^2} \frac{P_M}{D_2^2}. \quad (37)$$

The partial derivative of L , D_1 , and D_2 with regard to M are calculated as follows:

$$\frac{\partial L}{\partial M} = L(\rho - \ln \rho - 1) + LM(1 - \frac{1}{\rho}) \frac{\partial \rho}{\partial M},$$

$$\begin{aligned} \frac{\partial D_1}{\partial M} &= \sqrt{2\pi} \left(\frac{1}{2\sqrt{M}}(1-\rho)L + \sqrt{M} \left(-\frac{\partial \rho}{\partial M} \right) L + \sqrt{M}(1-\rho) \frac{\partial L}{\partial M} \right) \\ &= \sqrt{2\pi} \left(\frac{1}{2\sqrt{M}}(1+\rho)L - \sqrt{M}(1-\rho) \ln \rho L \right), \end{aligned}$$

$$\frac{\partial D_2}{\partial M} = -\frac{\partial \rho}{\partial M} = \frac{\rho}{M}.$$

Substitute $\frac{\partial L}{\partial M}$, $\frac{\partial D_1}{\partial M}$, and $\frac{\partial D_2}{\partial M}$ back into the Equation (37), we have

$$\frac{\partial g_1(M, s)}{\partial M} = \frac{\bar{x}_1}{M D_1 D_2} \left[\frac{\partial D_1}{\partial M} D_2 + \frac{\partial D_2}{\partial M} D_1 + \frac{1}{M} \right].$$

The partial derivative of $g_2(M, s)$ with regard to M can be easily calculated as

$$\frac{\partial g_2(M, s)}{\partial M} = (\xi s^\gamma - P_{sta}) \cdot \rho + P_{sta}.$$

Calculate the partial derivative of Equation (31) with regard to s : The partial derivative of L , D_1 , and D_2 with regard to s are calculated as follows:

$$\frac{\partial L}{\partial s} = LM(1 - \frac{1}{\rho}) \frac{\partial \rho}{\partial s} = \frac{LM}{s} (1 - \rho),$$

$$\frac{\partial D_1}{\partial s} = \sqrt{2\pi M} \left[\left(-\frac{\partial \rho}{\partial s} L + (1-\rho) \frac{\partial L}{\partial s} \right) \right] = \sqrt{2\pi M} [\rho + M(1-\rho)^2] \frac{L}{s},$$

$$\frac{\partial D_2}{\partial s} = -\frac{\partial \rho}{\partial s} = \frac{\rho}{s}.$$

The partial derivatives of $g_1(M, s)$ and $g_2(M, s)$ with regard to s are hence computed as

$$\frac{\partial g_1(M, s)}{\partial s} = -\frac{\bar{x}_1}{M} \cdot \frac{\partial}{\partial s} \left(\frac{P_M}{(1-\rho)^2} \right) = -\frac{\bar{x}_1}{M} \left[\frac{\frac{\partial D_1}{\partial s} D_2 + \frac{\partial D_2}{\partial s} D_1}{(D_1 D_2)^2} \right],$$

$$\frac{\partial g_2(M, s)}{\partial s} = M \rho \xi \gamma s^{\gamma-1}.$$

Once we obtain the above partial derivatives of Equation (31) with regard to ω , M and s , we can compute and obtain the optimal solutions by letting these partial derivatives of Equation (31) with regard to ω , M , and s equal 0.

4.2 Solve Augmented Lagrange Function

We present in this section an augmented Lagrange multiplier method based algorithm that solves the profit optimization problem given in (32) and derives its optimum solutions, including the service price and multiserver configurations. The proposed algorithm first computes an optimum Lagrange multiplier, which guarantees that the solution of original objective function and the solution of Lagrange function are consistent in the case where the optimal multiplier is obtained. Subsequently, the optimal service price and multiserver configurations are determined.

Let $M^{(k)}$, $s^{(k)}$, and $v^{(k)}$ indicate the k^{th} iteration of M , s , and v in the algorithm. Let ε , η , and Ψ be three positive numbers, l be the number of iterations, and L be the maximum number of iterations. Algorithm 1 describes the proposed augmented Lagrange algorithm. Inputs to the algorithm are electricity price C^τ during time slot τ , the rent δ , and user requests arrival rate λ_u . The algorithm iteratively derives the optimal cloud service price ω and multiserver configurations which includes the optimal number of servers M , the server speed s , and the *Profit* of the cloud service provider.

The algorithm works as follows. It first formulates the optimization problem into the form in Equation (25), then sets parameters of ε , η , Ψ , and L , and initializes variables of $M^{(0)}$, $s^{(0)}$, v^1 , and l (lines 1-3). In each round of iteration, the algorithm calls the augmented Lagrange function solver, denoted by **ALF-Solver**($\phi(\omega, M^{(l-1)}, s^{(l-1)}, v^{(l)}, \sigma)$), to obtain a local optimum of the ω , M , and s (line 5). The **ALF-Solver**($\phi(\omega, M^{(l-1)}, s^{(l-1)}, v^{(l)}, \sigma)$) derives the local optimum by computing partial derivatives of $\phi(\omega, M, s, v, \sigma)$ with regard to ω , M , and s , and solving a system of equations of ω , M , and s (lines 18-21).

The algorithm exits if the Lagrange multiplier vector \mathbf{v} converges and approximates the optimum by an error of ε . Let $Q_j(M^{(l)}, s^{(l)}) = g_j(M^{(l)}, s^{(l)}) - y_j^2$ for $j = 1, 2$ be the penalty item of the augmented Lagrange function given in Equation (32), then the Lagrange multiplier vector \mathbf{v} converges if $\|Q(M^{(l)}, s^{(l)})\| < \varepsilon$ holds (lines 6-10). If it does not converge or converges too slowly, that is, $\|Q(M^{(l)}, s^{(l)})\| / \|Q(M^{(l-1)}, s^{(l-1)})\| \geq \Psi$ holds for a positive number Ψ , the penalty factor σ is updated to $\eta\sigma$ for $\eta > 1$ to speed up the convergence process (lines 11-13). Accordingly, the Lagrange multiplier for the next iteration is updated to $v_j^{l+1} = \max(0, v_j^{(l)} - \sigma g_j(M^{(l)}, s^{(l)}))$ ($j = 1, 2$) (lines 14-15), and the procedure moves to the next iteration.

Algorithm 1: Iteratively solve the augmented Lagrange function

Input:
Electricity price C^τ during sales period τ , rent δ , user requests arrival rate λ_u ;

Output:
The optimal service price ω , number of servers M , server speed s , and *Profit*;

- 1 Formulate the optimization problem into the form in Equation (25);
- 2 Set parameters $\alpha, \beta, \gamma, \varepsilon, \eta, \Psi$, and L ;
- 3 Initialize $M^{(0)}$, $s^{(0)}$, $v^{(1)}$, and $l = 1$;
- 4 **while** $l < L$ **do**
- 5 $[\omega^{(l)}, M^{(l)}, s^{(l)}] = \mathbf{ALF-Solver}(\phi(\omega, M^{(l-1)}, s^{(l-1)}, v^{(l)}, \sigma))$;
 // exit when $\{v^{(l)}\}$ converges;
- 6 **if** $\|Q(M^{(l)}, s^{(l)})\| < \varepsilon$ **then**
- 7 Calculate the *Profit* using the Equation (24);
 // record the optimal solution;
- 8 $Result[\theta_{sub}] = [Profit, \omega^{(l)}, M^{(l)}, s^{(l)}]$;
- 9 **break**;
- 10 **end**
 // otherwise, increase penalty factor σ ;
- 11 **else if** $\|Q(M^{(l)}, s^{(l)})\| / \|Q(M^{(l-1)}, s^{(l-1)})\| \geq \Psi$ **then**
- 12 $\sigma = \eta\sigma$;
- 13 **end**
 // update the multiplier vector \mathbf{v} ;
- 14 $v_j^{l+1} = \max(0, v_j^{(l)} - \sigma g_j(M^{(l)}, s^{(l)}))$ ($j = 1, 2$);
- 15 $l = l + 1$;
- 16 **end**
- 17 **return** $Result[\theta_{sub}]$;
 // solve the Lagrange function in (35);
- 18 **ALF-Solver**($\phi(\omega, M^{(l-1)}, s^{(l-1)}, v^{(l)}, \sigma)$)
- 19 Compute partial derivatives of ϕ w.r.t. ω, M , and s as $\partial\phi(\omega, M^{(l-1)}, s^{(l-1)}, v^{(l)}, \sigma) / \partial(\omega, M, s)$;
- 20 Calculate ω, M , and s based on a system of equations of $\frac{\partial\phi}{\partial\omega}, \frac{\partial\phi}{\partial M}$, and $\frac{\partial\phi}{\partial s}$;
- 21 **return** $[\omega, M, s]$;

Once the algorithm converges, the optimum of ω , M , and s are derived, and the optimal *Profit* of the cloud service provider can be calculated by using Equation (24) (line 7). Line 17 returns the optimal service price, multiserver configurations, and *Profit* of the cloud service provider.

4.3 Design a Dynamic Closed Loop Control Scheme

The solution to the profit maximization problem described above focuses on the interaction between users and the cloud service provider. However, the impact of dynamic cloud computing environment such as fluctuating electricity bill and rental fees on profit maximization mechanism is not investigated. On one hand, the variation of electricity bill or rental fees has a direct impact on the expenditure of the cloud service provider. On the other hand, the variation of electricity bill or rental fees has an indirect influence on user perceived value which affects the user demand of the cloud service, and ultimately impacts the revenue of the cloud service provider. Thus, it is necessary to design a scheme to adjust service price and multiserver configurations according to the dynamics of the cloud computing environment.

In this subsection, we propose a closed loop control scheme to dynamically update the optimal service price and multiserver configurations. As illustrated in Figure 2, the runtime control scheme monitors the dynamic cloud computing environment. Once the electricity bill or rental fees changes, the proposed control scheme first fits a new probability distribution function of user perceived value using kernel density estimation based on the historical price data set Ω . Subsequently, it reconstructs and resolves the profit maximization problem based on the new probability distribution and the variation of electricity bill or rental fees.

The kernel density estimation technique adopted in the proposed control scheme is a stochastic non-parametric way to estimate the probability density function of a random variable [26]. It is a fundamental data smoothing technique where inferences about the population are made based on a finite data sample. Given a univariate independent and identically distributed sample drawn from some distribution with an unknown density function, the technique can be used to estimate the shape of the density function.

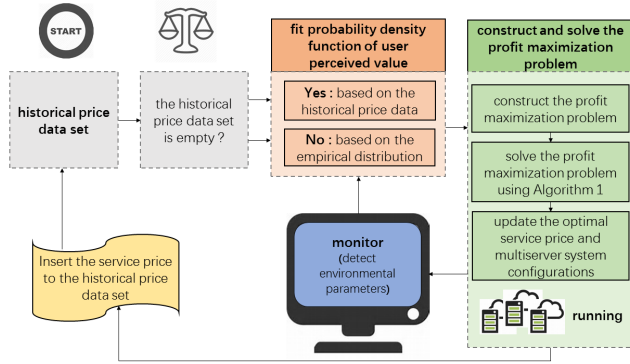


Figure 2: Overview of closed loop control scheme

The details of the proposed runtime control scheme are described in Algorithm 2. Inputs to the algorithm are the historical price data set Ω and the output of system monitor. The closed loop control scheme works as follows. It monitors whether parameters of the cloud computing environment change at all times (line 2). If no change, the system will run with the current multiserver configurations (lines 3-5). Otherwise, it updates and solves the profit maximization problem (lines 6-15). Lines 7-8 fit the probability density function (pdf) of user perceived value using MATLAB function `ksdensity(*)` based on historical price data set Ω . Line 9 updates the profit maximization problem according to the change of the cloud computing environment (i.e., electricity bill or rental fees). Line 10 solves the profit maximization problem using algorithm 1. The algorithm updates the optimal cloud service price and multiserver configurations in line 11, and calculates the profit of the cloud service provider using Equation (24) in line 12. Finally, it inserts the service price ω into the historical price data set Ω in line 13. Line 14 returns the optimal cloud service price, multiserver configurations, and *Profit* of the cloud service provider.

The MATLAB function `ksdensity(*)` is used to fit the probability density function of user perceived value based on historical price data by using kernel smoothing density

estimation. Line 18 first calculates the number of samples in the historical price data set Ω . The kernel bandwidth h is a free parameter that exhibits a strong influence on the resulting estimate [26]. Here, Gaussian basis functions are used to approximate univariate data. Thus, the optimal choice for kernel bandwidth h is calculated as line 19, which minimizes the mean integrated squared error used in density estimation [43]. `std(Ω)` computes the standard deviation of the samples in Ω . Line 20 uses operator `@(x)` to define the function handle *phi*, which represents a normal probability density function. `exp(x)` and `sqrt(x)` represent exponential function and square root function, respectively. Based on normal probability density function *phi* and bandwidth h , line 21 computes the kernel density, that is probability density function by defining the function handle *ksden*. `mean(x)` is used to compute the average of the array. Line 22 returns the final fitted probability density function.

Algorithm 2: Dynamic closed loop control scheme

Input:

The historical price data set Ω , the output of system monitor;

Output:

The optimal service price ω , number of servers M , server speed s , and *Profit*;

```

1 while true do
2   Monitor if parameters of computing environment
   change;
3   if no change then
4     continue;
5   end
6   else
7      $\Omega$  = historical price data set;
   // fit pdf using ksdensity( $\Omega$ );
8      $f_X(\omega) \leftarrow \text{ksdensity}(\Omega)$ ;
9     Update profit maximization problem given in (25);
10    Solve profit maximization problem using Algo-
   rithm 1;
11    Update the optimal service price  $\omega$ , number of
   servers  $M$ , and server speed  $s$ ;
12    Calculate the Profit using Equation (24);
13    Insert service price  $\omega$  into historical price data set
    $\Omega$ ;
14    return [ $\omega, M, s, Profit$ ];
15  end
16 end
   // fit pdf of user perceived value using
   MATLAB function ksdensity( $\Omega$ );
17 ksdensity( $\Omega$ )
   // derive the number of samples in  $\Omega$ ;
18  $n = \text{length}(\Omega)$ ;
   // set the optimal bandwidth  $h$ ;
19  $h = \text{std}(\Omega) * (4/3/n)^{(1/5)}$ ;
   // obtain the normal pdf;
20  $phi = @(x)(\exp(-.5 * x.^2)/\text{sqrt}(2 * pi))$ ;
   // compute kernel density with phi and  $h$ ;
21  $ksden = @(x)\text{mean}(phi((x - \Omega)/h)/h)$ ;
22 return ksden;

```

5 SIMULATION-BASED EVALUATION

Extensive simulation experiments have been conducted to validate the effectiveness of the proposed scheme. We

first describe simulation settings in detail, then verify the effectiveness of the proposed user perceived value-based dynamic pricing model, followed by the validation of the optimal pricing and multiserver configurations and a comparison study with benchmarking schemes in terms of the profit of the cloud service provider.

5.1 Simulation Settings

The simulation experiments are conducted on a machine equipped with 2.56GHz Intel i7 quad-core processor and 8GB DDR4 memory, and running a Windows version of Matlab_x64. For the sake of a fair comparison, three types of users used in [35] are also adopted in our simulation experiments. Users of type 1 are delay-sensitive while users of type 2 and 3 are delay-insensitive to the deferment of the service requests. Data of type 1 were extracted from Youtube U.S. traffic from January 1, 2014 to January 31, 2014 [44]. Data of type 2 and 3 were extracted from GMaps and Gmail U.S. traffic from January 1, 2014 to January 31, 2014 [44], respectively. The one day ahead real-time pricing data released by Ameren Illinois Power Corporation at January 2014 are taken as the price input in the experiment [45]. We also assume that user perceived value X obeys the following normal distribution, $X \sim N(0, 0.22)$ [7], [22]. In addition, the value of other parameters used in our simulation experiment are shown in Table 1.

Table 1: Experimental parameters table

Parameter	Definition	Value
T	sales period	30 d
τ	time slot	1 h
N	number of time slot	720
D_{max}	maximum value of service deferment	24
ψ_1	sensitivity factor of users of type 1	∞
ψ_2	sensitivity factor of users of type 2	0.1
ψ_3	sensitivity factor of users of type 3	0.11

5.2 Verify User Perceived Value-Based Dynamic Pricing Model

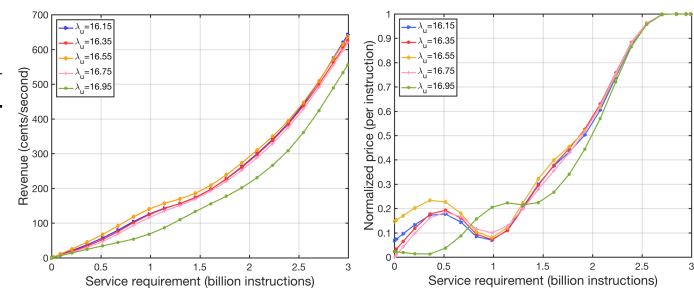
This subsection verifies the proposed user perceived value-based dynamic pricing model from the perspective of supply and demand law.

5.2.1 Revenue Vs. Service Requirement

We first analyze the relationship between the service requirement in terms of the number of instructions, which is denoted by r , and the revenue of the cloud service provider. In addition to the parameters given in Table 1, we set the average service requirement denoted by \bar{r} to 1 billion instructions. The number of servers M is initialized to 7, the base speed s_0 and speed s of servers are both initialized to 1 billion instructions per second, and the static power consumption P_{sta} is set to 2W. The parameters of dynamic power consumption are assumed to be $\gamma = 2.0$ and $\xi = 9.4192$, and parameters of Gamma distribution are assumed to be $\alpha = 2.0$ and $\beta = 1.5$ [6].

Figure 3(a) shows the relationship between service requirements and the revenue of the cloud service provider when service request arrival rate λ_u is 16.15, 16.35, 16.55, 16.75, and 16.95 billions instructions per second, respectively. It can be seen from Figure 3(a) that the revenue increases as service requirements increase. This indicates that the usage of cloud services and the revenue obtained are positively correlated under the user perceived value-based pricing model. In addition, as shown in the figure, the revenue decreases as λ_u increases. This is because with the increase of λ_u , servers can not process service requests in time, leading to a higher response time and lower quality of service. Low quality of service will result in a smaller number of users to purchase the cloud services, thus, the revenue of the cloud service provider decreases accordingly.

Figure 3(b) shows the relationship between service requirements and the normalized service price when service request arrival rate λ_u is 16.15, 16.35, 16.55, 16.75, and 16.95 billions instructions per second, respectively. From the figure, we can see that when service requirement $r < 1.4$ billions, the normalized price fluctuates with the increase of the service requirement. When service requirement $1.4 < r < 2.6$ billions, the normalized price increases with the increase of the service requirement. When service requirement $r > 2.6$ billions, the normalized price eventually converges to a stable value with the increase of the service requirement.



(a) Revenue vs. service requirement. (b) Normalized price vs. service requirement.

Figure 3: Relationship between service requirement and revenue/normalized price.

5.2.2 Purchase Amount and Revenue Vs. Service Price

Figure 4(a)-4(d) demonstrate that how the relationship among the cloud service purchase amount, revenue, and the price of cloud service changes when service request arrival rate λ_u is 16.15, 16.55, 16.75, and 16.95 billions instructions per second, respectively. As we can see from these figures, before the service price reaches user perceived value of the service, the purchase amount of the cloud service increases with the increases of the price. Once the price exceeds user perceived value of the service, the purchase amount declines sharply. This observation is consistent with real market situation, that is, users are willing to accept a price and purchase when the price is lower than their perceived value. However, the user's purchase intention will decline sharply when the price is beyond user perceived value.

It also can be seen from Figure 4(a) to 4(d) that the point where purchase amount is maximum is not necessarily the point where the revenue is maximum. That is, the revenue

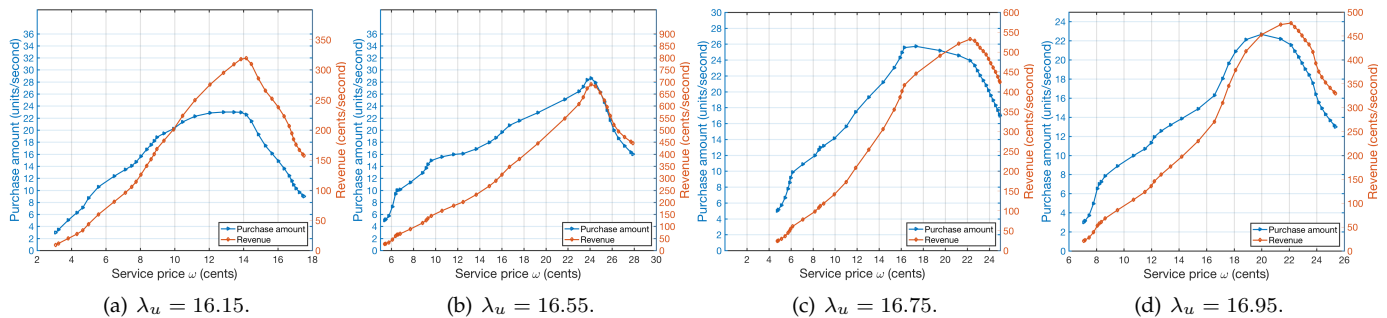
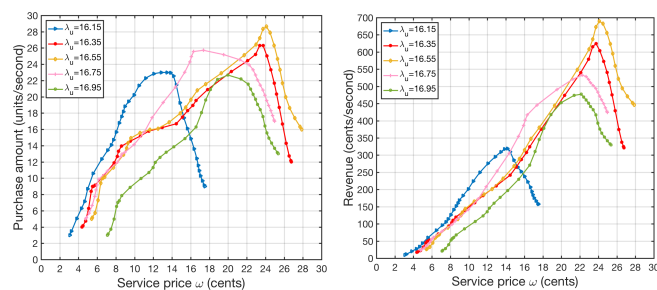


Figure 4: Purchase amount and revenue vs. service price.

for the scenario of the low price and high purchase amount is not necessarily higher than the revenue for the scenario of the high price and low purchase amount. When service request arrival rate $\lambda_u = 16.55$, the cloud service provider can get the maximum revenue.

5.2.3 Purchase Amount and Revenue Vs. Service request arrival rate

Figure 5(a) and 5(b) demonstrate that how the cloud service purchase amount and revenue change when service request arrival rate λ_u is 16.15, 16.35, 16.55, 16.75, and 16.95 billions instructions per second, respectively. From Figure 5(a), we can see that the optimal prices for the maximum service purchase are different under diverse service request arrival rate λ_u . For the case where $\lambda_u = 16.55$ billions instructions per second, the service purchase amount and the service price ω reach the maximum value at the same time when compared to cases of different service request arrival rates. Meanwhile, the maximum purchase amount at $\lambda_u = 16.35$ is approximately the same as the maximum purchase amount at $\lambda_u = 16.75$. This situation holds for the case where $\lambda_u = 16.15$ and $\lambda_u = 16.95$. This is because with the increase of λ_u , limited number of servers can not process arrived service requests in time, leading to a higher response time, lower quality of service, and thus a lower maximum purchase amount of cloud services.



(a) Purchase amount vs. service price (b) Revenue vs. service request arrival rate.

Figure 5: Purchase amount and revenue vs. service request arrival rate.

The revenue in Figure 5(b) is obtained by multiplying the purchase amount and service price in Figure 5(a). Figure 5(b) shows that the optimal prices for the maximum revenue are different under various service request arrival rate λ_u .

From this figure, we observe that with the increase of λ_u , the maximum revenue at different λ_u increases first and then decreases. The cloud service provider obtains the maximum revenue when $\lambda_u = 16.55$ billions instructions per second. Similarly, this is because with the increase of λ_u , limited number of servers can not process arrived service requests in time, leading to a lower maximum purchase amount of cloud services, and thus a lower revenue. Based on the above experimental results, our user perceived value-based dynamic pricing model follows the supply and demand law in market.

5.3 Validate Multiserver Configurations for Profit Maximization

We set the response time constraint for user requests, denoted by b_1 , to 0.33 seconds and the power consumption of the server system, denoted by b_2 , to 10^6 W. The rental cost denoted by δ is set to 1.5 cents per second [7].

Figure 6(a) shows the relationship between profit and the number of working servers. It can be seen from the figure that when service request arrival rate $\lambda_u = 12.9, 13.9, 14.9, 15.9,$ and 16.9 billions instructions per second, the optimal number of servers denoted by M is 16, 17, 19, 18, and 17, respectively. It is clear that when M is small, the utilization of working servers is approaching 1, leading to a long response time for user requests and low quality of service accordingly, and in turn a low profit under the user perceived value-based dynamic pricing model. As M increases, the number of user requests in the waiting queue decreases quickly, the user requests do not have to wait too long, and thus the profit increases under the user perceived value-based dynamic pricing model. However, as M continues increasing, the profit does not increase. This is because the increase in the number of servers leads to an increase in the maintenance cost of working servers including electricity and rental cost.

Figure 6(b) shows the relationship between profit and the server speed s . We notice from the figure that in order to maximize the profit, the optimal speed s is set to 0.7642, 0.9435, 1.1044, 1.1293, and 1.2838 billions instructions per second when the service request arrival rate $\lambda_u = 12.9, 13.9, 14.9, 15.9,$ and 16.9 billions instructions per second, respectively. It is clear that when the server speed s is low, the utilization of working servers is approaching 1, leading to a long response time for user requests and low quality of service accordingly, and in turn a low profit under the user

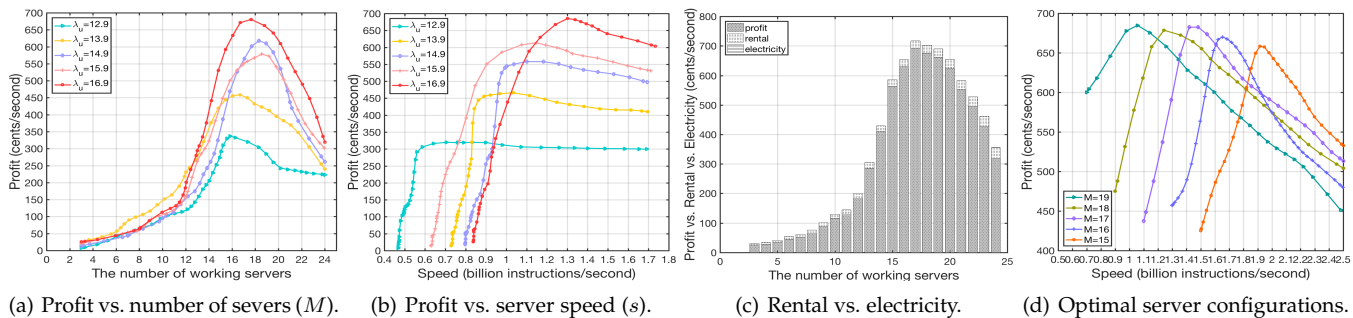


Figure 6: Validate server configurations for profit maximization.

perceived value-based dynamic pricing model. When the server speed s is high, service requests are more likely to be executed on time, leading to an increase in the profit under the user perceived value-based dynamic pricing model. However, with the continued increase in s , the profit does not increase as expected. This is because the increase in the server speed leads to an increase in the cost of operating a multiserver system.

From Figure 6(a) and 6(b), we see that the profit reaches its maximum when λ_u is 16.9 billions instructions per second and the number and speed of servers take the appropriate value. In addition, we observe such a phenomenon that the profit in Figure 6(a) drops faster than in Figure 6(b) after reaching the maximum value. From Section 2, we know that the number of servers M affects both rental fees and electricity bill while the server speed s only impacts electricity bill. Figure 6(c) studies the impact of rental fees and electricity bill on profit under the same experimental conditions. We can see from the figure that the rental fees have a greater impact on profit compared to electricity bill. Thus, with the increase of M , the profit in Figure 6(a) decreases faster due to the great impact of rental fees. Meanwhile, with the increase of s , the profit in Figure 6(b) decreases slower due to the weak impact of electricity bill.

Figure 6(d) gives the optimal M and s of servers that maximize the profit when $\lambda_u = 16.9$ billions instructions per second. It can be seen that the maximal profit is obtained when s and M is set to 1.4351 billions instructions per second and 17, respectively. That is to say, 687.9 cents of profit is obtained when 17 servers are open and each server runs at 1.4351 billions instructions per second.

5.4 Compare the Maximal Profit with Benchmarking Pricing Strategies

We compare the proposed user perceived value-based profit maximization scheme with two benchmarking methods OMCPM [6] and UPMR [35]. OMCPM [6] is an efficient pricing model that takes such factors into considerations as the service-level agreement and customer satisfaction. It derives an optimal server configuration and service price for profit maximization. UPMR [35] is a usage based pricing model used by today's major cloud service providers. The UPMR model rewards users proportionally based on the time length that users set as deadlines for completing their service requests. Compared with OMCPM and UPMR, our

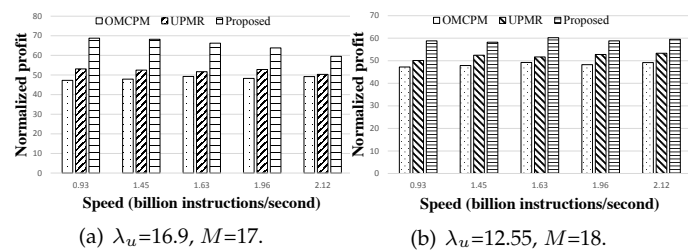


Figure 7: Compare the maximal profit with two benchmarking pricing models.

pricing method is based on user perceived value that reflects users willingness to purchase cloud services.

We compare the maximal profit generated by proposed pricing model with that generated by the two benchmarking pricing models under the same experimental settings. Two comparison experiments are conducted. In the first experiment, user service request arrival rate λ_u is set to 16.9 billions instructions per second and the number of working servers M is set to 17. In the second experiment, λ_u is set to 12.55 billions instructions per second and M is set to 18. It is clear from Figure 7 that our proposed dynamic pricing model is superior to the two benchmarking models. For instance, the proposed pricing model can obtain up to 21.55 cents per second more (31.32%) as compared to OMCPM method, and 15.66 cents per second more (22.76%) as compared to UPMR when $\lambda_u = 16.9$ billions instructions per second, $M = 17$, and $s = 0.93$ billion instructions per second. Thus, the pricing strategy based on user perceived value can better reflect the market demand and the cloud service provider can obtain higher profit.

We further verify how the expected number of actual buyers ($E_\omega(m)$) and the corresponding revenue change when user perceived value obeys normal distributions with different parameters. Figure 8(a)-8(f) show the expected number of actual buyers ($E_\omega(m)$) under different expectations μ and variances σ^2 of user perceived value in our proposed dynamic pricing model. From Figure 8(a)-8(c), we can see that under different expectations μ , as μ increases, the cloud service provider needs to increase the service price ω to obtain the same amount of purchases. From Figure 8(d)-8(f), we can see that under different variances σ^2 , when service price ω is less than μ , as σ^2 increases, the cloud service provider needs to decrease the service

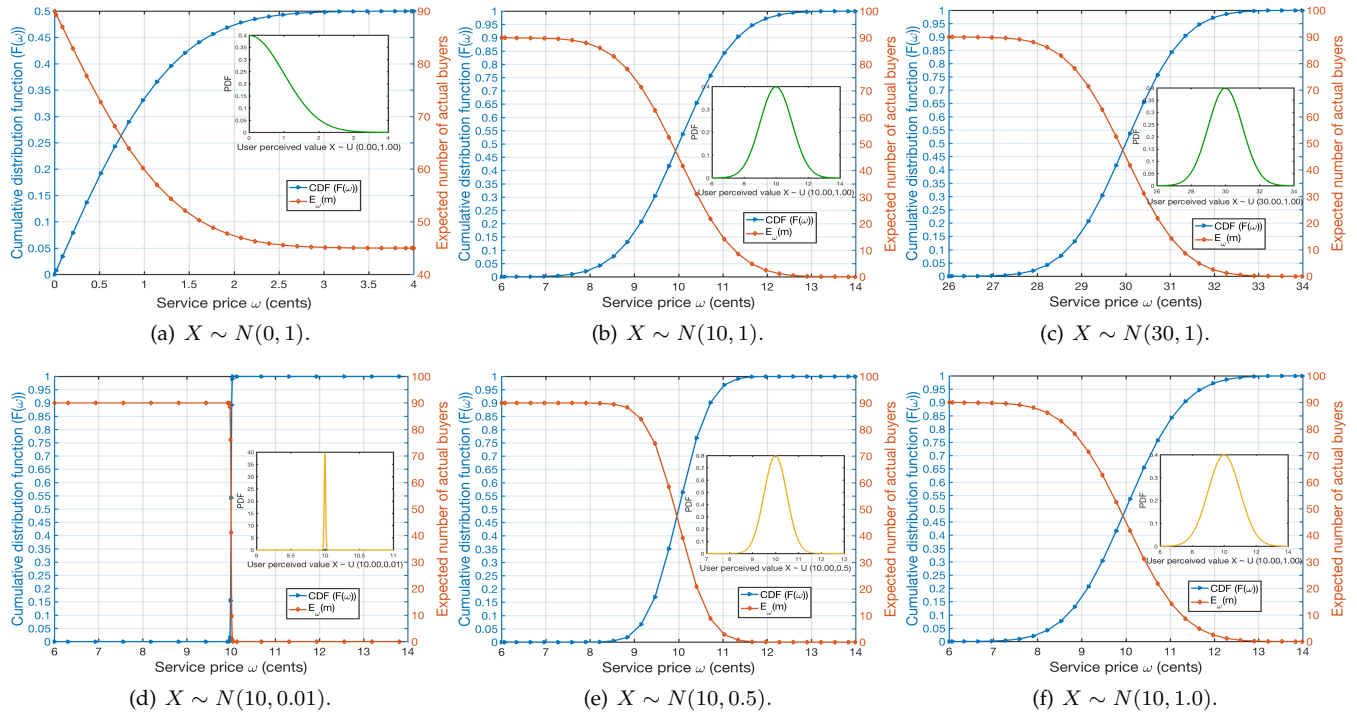


Figure 8: Verify the change in the expected number of actual buyers when user perceived value obeys normal distributions with different parameters (i.e., μ and σ^2).

price ω to obtain the same amount of purchases. However, when service price ω is greater than μ , as σ^2 increases, the cloud service provider needs to increase the service price ω to obtain the same amount of purchases. This is because the larger the σ^2 , the more dispersed the perceived value's distribution. Thus, in the case of the same service price ω , the purchase amount changes accordingly.

Figure 9(a)-9(f) show the revenue under different expectations μ and variances σ^2 of user perceived value in our proposed dynamic pricing model. From Figure 9(a)-9(c), we can find that under the same purchase amount, the cloud service provider needs to increase expectation μ of normal distribution, that is, users' perceived value of services, to achieve higher revenue. From Figure 9(d)-9(f), we can see that under the same purchase amount, the cloud service provider needs to decrease variance σ^2 of normal distribution to achieve higher revenue. In general, to obtain the higher revenue, the cloud service provider needs to carry out market strategies to improve perceived value of service in users' mind. This is because under the same purchase amount, that is, under the same number of requests that the cloud service provider needs to process, the corresponding expenses are the same. Thus, it is reasonable to grow the profit of the cloud service provider by increasing the revenue.

6 CONCLUSION

In this paper, we first propose a user perceived value-based dynamic profit maximization mechanism that takes into account the interaction between cloud users and the cloud service provider. Subsequently, we use augmented Lagrange multiplier method to solve the optimization problem to

derive the optimal solution, including the service price, number of servers, and speed of servers. Finally, we propose a dynamic closed loop control scheme to update the service price and multiserver configurations using kernel density estimation method. Extensive experimental results show that our proposed profit maximization scheme follows the supply and demand law in market, and are able to obtain more profit of up to 31.32% and 22.76% as compared to the state of the art benchmarking methods OMCPM [6] and UPMR [35], respectively.

REFERENCES

- [1] K. Hwang, J. Dongarra, and G. Fox, Distributed and cloud computing, *Morgan Kaufmann*, 2012.
- [2] L. Wang, H. Zhong, R. Ranjan, A. Zomaya, and P. Liu, Estimating the statistical characteristics of remote sensing big data in the wavelet transform domain, *IEEE Transactions on Emerging Topics in Computing*, vol. 2, no. 3, pp. 324-337, 2014.
- [3] P. Mell and T. Grance, The NIST definition of cloud computing, *Communications of the ACM*, vol. 15, 2011.
- [4] L. Wang, Y. Ma, A. Zomaya, R. Ranjan, and D. Chen, A parallel file system with application-aware data layout policies for massive remote sensing image processing in digital earth, *IEEE Transactions on Parallel and Distributed Systems*, vol. 26, no.6, pp. 1497-1508, 2015.
- [5] P. Cong, L. Li, G. Shao, J. Zhou, M. Chen, K. Huang, and T. Wei, User perceived value-aware data pricing for profit maximization of multiserver systems, *IEEE International Conference on Parallel and Distributed Systems*, pp. 537-544, 2017.
- [6] J. Cao, K. Hwang, K. Li, and A. Zomaya, Optimal multiserver configuration for profit maximization in cloud computing, *IEEE Transactions on Parallel and Distributed Systems*, vol. 24, no. 6, pp. 1087-1096, 2013.
- [7] Y. Chun, Optimal pricing and ordering policies for perishable commodities, *European Journal of Operational Research*, pp. 68-82, 2003.
- [8] C. Li, Cloud computing system management under flat rate pricing, *Journal of network and systems management*, vol. 19, no. 3, pp. 305-318, 2011.

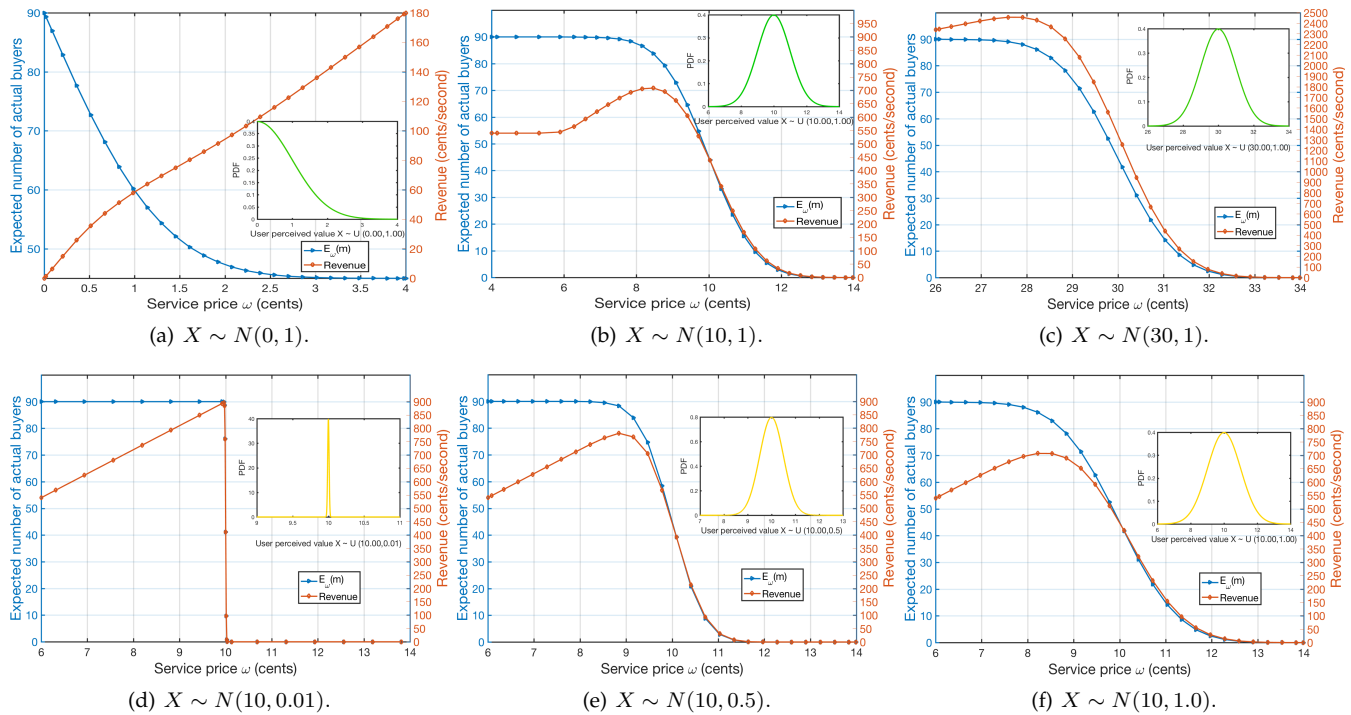


Figure 9: Verify the change in the revenue when user perceived value obeys normal distributions with different parameters (i.e., μ and σ^2).

[9] G. Kesidis, A. Das, and G. Veciana, On flat-rate and usage-based pricing for tiered commodity internet services, *Annual Conference on Information Sciences and Systems*, pp. 304-308, 2008.

[10] Y. Lee, C. Wang, A. Zomaya, and B. Zhou, Profit-driven scheduling for cloud services with data access awareness, *Journal of Parallel and Distributed Computing*, vol. 72, no. 4, pp. 591-602, 2012.

[11] M. Macias and J. Guitart, A genetic model for pricing in cloud computing markets, *ACM Symposium on Applied Computing*, pp. 113-118, 2011.

[12] Amazon EC2. [Online]. Available: <http://aws.amazon.com>.

[13] Amazon EC2 spot instances. [Online]. Available: <https://aws.amazon.com/cn/ec2/spot/pricing>.

[14] H. Xu and B. Li, Dynamic cloud pricing for revenue maximization, *IEEE Transactions on Cloud Computing*, vol. 1, no. 2, pp. 158-171, 2013.

[15] J. Zhao, H. Li, C. Wu, Z. Li, Z. Zhang, and F. Lau, Dynamic pricing and profit maximization for the cloud with geo-distributed data centers, *IEEE Conference on Computer Communications*, 2014.

[16] Service-level agreement. [Online]. Available: https://en.wikipedia.org/wiki/Service-level_agreement.

[17] M. Ghamkhari and H. Mohsenian-Rad, Energy and performance management of green data centers: A profit maximization approach, *IEEE Transactions on Smart Grid*, vol. 4, no. 2, pp. 1017-1025, 2013.

[18] Y. Lee, C. Wang, A. Zomaya, and B. Zhou, Profit-driven service request scheduling in clouds, *International Conference on Cluster, Cloud and Grid Computing*, pp. 15-24, 2010.

[19] D. Irwin, L. Grit, and J. Chase, Balancing risk and reward in a market-based task service, *International Conference on high performance distributed computing*, pp. 160-169, 2004.

[20] L. Wang, Y. Ma, J. Yan, V. Chang, A. Zomaya, pipsCloud: High performance cloud computing for remote sensing big data management and processing, *Future Generation Computer Systems*, pp. 353-368, 2018.

[21] M. Leppaniemi, H. Karjaluoto, and H. Saarijarvi, Customer perceived value, satisfaction, and loyalty: the role of willingness to share information, *The International Review of Retail, Distribution and Consumer Research*, vol. 27, no. 2, pp. 164-188, 2017.

[22] Z. Yang and R. Peterson, Customer perceived value, satisfaction and loyalty: The role of switching costs, *Psychology & Marketing*, vol. 21 no. 10, pp. 799-822, 2004.

[23] S. Karlin and C. Carr, Prices and optimal inventory policy, *Arrow Karlin & Scarf Studies in Applied Probability & Management Science*, pp. 159-172, 1962.

[24] J. Roig, J. Garcia, M. Tena, and J. Monzonis, Customer perceived value in banking services, *International Journal of Bank Marketing*, vol. 24, no. 5, pp. 266-283, 2006.

[25] Definition of Perceived Value. [Online]. Available: <http://smallbusiness.chron.com/definition-perceived-value-23017.html>.

[26] Kernel density estimation. [Online]. Available: https://en.wikipedia.org/wiki/Kernel_density_estimation.

[27] J. Gallego and G. Ryzin, Optimal dynamic pricing of inventories with stochastic demand over finite horizons, *Management Science*, vol. 40, no. 8, pp. 999-1020, 1994.

[28] P. Pfeiffer. Probability for applications, Springer, 2012.

[29] K. Li, Optimal load distribution for multiple heterogeneous blade servers in a cloud computing environment, *Journal of Grid Computing*, pp. 943-952, 2011.

[30] B. Chun and D. Culler, User-centric performance analysis of market-based cluster batch schedulers, *International Symposium on CLUSTER Computing and the Grid*, 2002.

[31] K. Li, Optimal configuration of a multicore server processor for managing the power and performance tradeoff, *Journal of Supercomputing*, vol. 61, no. 1, pp. 189-214, 2012.

[32] L. Kleinrock, Queueing systems, Volume 1: Theory, Wiley, 1975.

[33] J. Zhou, J. Chen, K. Cao, T. Wei, and M. Chen, Game theoretic energy allocation for renewable powered in-situ server systems, *IEEE International Conference on Parallel and Distributed Systems*, pp. 721-728, 2016.

[34] J. Little and S. Graves, Little's law, *International Series in Operations Research and Management Science*, pp. 81-100, 2008.

[35] Y. Zhan, M. Ghamkhari, D. Xu, S. Ren, and H. Mohsenian-Rad, Extending demand response to tenants in cloud data centers via non-intrusive workload flexibility pricing, *IEEE Transactions on Smart Grid*, pp. 1-8, 2016.

[36] Steepest descent method - wikipedia. [Online]. Available: https://en.wikipedia.org/wiki/Method_of_steepest_descent.

[37] Newton's method - wikipedia. [Online]. Available: https://en.wikipedia.org/wiki/Newton%27s_method.

[38] D. Bertsekas, Multiplier methods: A survey, *Automatica*, vol. 12, pp. 133-145, 1976.

- [39] W. Long, X. Liang, S. Cai, J. Jiao, and W. Zhang, A modified augmented Lagrangian with improved grey wolf optimization to constrained optimization problems, *Neural Computing & Applications*, pp. 1-18, 2016.
- [40] Y. Zheng and Z. Meng, A new augmented Lagrangian objective penalty function for constrained optimization problems, *Open Journal of Optimization*, vol. 6, pp. 39-46, 2017.
- [41] B. Dandurand, N. Boland, J. Christiansen, A. Eberhard, and F. Oliveira, A parallelizable augmented Lagrangian method applied to large-scale non-convex-constrained optimization problems, 2017.
- [42] Completing the square-wikipedia. [Online]. Available: https://en.wikipedia.org/wiki/Completing_the_square.
- [43] Mean integrated squared error. [Online]. Available: https://en.wikipedia.org/wiki/Mean_integrated_squared_error.
- [44] Browse real-time traffic to google products and services. [Online]. Available: <http://www.google.com/transparencyreport/traffic/explorer>.
- [45] Real time prices-america. [Online]. Available: <https://www.america.com/RetailEnergy/RealTimePrices>.



Peijin Cong received her B.S. degree from the Department of Computer Science and Technology, East China Normal University, Shanghai, China, in 2016. She is currently pursuing the master degree with the Department of Computer Science and Technology, East China Normal University, Shanghai, China. Her current research interest is in the areas of power management in mobile devices and edge computing.



Liying Li received her B.S. degree from the Department of Computer Science and Technology, East China Normal University, Shanghai, China, in 2017. She is currently pursuing the master degree with the Department of Computer Science and Technology, East China Normal University, Shanghai, China. Her current research interests are in the areas of cyber physical systems and IoT resource management.



Junlong Zhou received his Ph.D. degree in Computer Science from East China Normal University, Shanghai, China, in 2017. He was a Research Visitor with the University of Notre Dame, Notre Dame, IN, USA, during 2014-2015. He is currently an Assistant Professor with the School of Computer Science and Engineering, Nanjing University of Science and Technology, Nanjing, China. His research interests include real-time embedded systems, cyber physical systems, and cloud computing. Dr. Zhou has been an

Associate Editor for the Journal of Circuits, Systems, and Computers since 2017. He is a member of the IEEE.



Kun Cao is currently pursuing his Ph.D. degree with the Department of Computer Science and Technology, East China Normal University, Shanghai, China. His current research interests are in the areas of high performance computing, multiprocessor systems-on-chip and cyber physical systems.



Tongquan Wei received his Ph.D. degree in Electrical Engineering from Michigan Technological University in 2009. He is currently an Associate Professor in the Department of Computer Science and Technology at the East China Normal University. His research interests are in the areas of Internet of Things, real-time embedded systems, green and reliable computing, parallel and distributed systems, and cloud computing. He serves as a Regional Editor for Journal of Circuits, Systems, and Computers since 2012.

He is a member of the IEEE.



Mingsong Chen (S'08-M'11) received the B.S. and M.E. degrees from Department of Computer Science and Technology, Nanjing University, Nanjing, China, in 2003 and 2006 respectively, and the Ph.D. degree in Computer Engineering from the University of Florida, Gainesville, in 2010. He is currently a full Professor with the Department of Embedded Software and Systems of East China Normal University. His research interests are in the area of design automation of cyber-physical systems, formal verification techniques and mobile cloud computing. He is a member of the IEEE.

He is a member of the IEEE.



Shiyan Hu received his Ph.D. in Computer Engineering from Texas A&M University in 2008. He is an Associate Professor at Michigan Tech, and he was a Visiting Associate Professor at Stanford University from 2015 to 2016. His research interests include Cyber-Physical Systems (CPS), CPS Security, Data Analytics, and Computer-Aided Design of VLSI Circuits, where he has published more than 100 refereed papers. He is an ACM Distinguished Speaker, an IEEE Systems Council Distinguished Lecturer,

an IEEE Computer Society Distinguished Visitor, and a recipient of National Science Foundation (NSF) CAREER Award. Prof. Hu is the Chair for IEEE Technical Committee on Cyber-Physical Systems. He is the Editor-In-Chief of IET Cyber-Physical Systems: Theory & Applications. He is an Associate Editor for IEEE Transactions on Computer-Aided Design, IEEE Transactions on Industrial Informatics, and IEEE Transactions on Circuits and Systems. He is also a Guest Editor for a number of IEEE/ACM Journals such as Proceedings of the IEEE and IEEE Transactions on Computers. He has held chair positions in numerous IEEE/ACM conferences. He is a Fellow of IET.