

UTILITY OPTIMIZATION-BASED BANDWIDTH ALLOCATION FOR ELASTIC AND INELASTIC SERVICES IN PEER-TO-PEER NETWORKS

SHIYONG LI^a, YUE ZHANG^a, YAN WANG^a, WEI SUN^{a,*}

^aSchool of Economics and Management
Yanshan University, No. 138 Hebei Avenue, Qinhuangdao 066004, China
e-mail: wsun@ysu.edu.cn

This paper considers reasonable bandwidth allocation for multiclass services in peer-to-peer (P2P) networks, measures the satisfaction of each peer as a customer by a utility function when acquiring one service, and develops an optimization model for bandwidth allocation with the objective of utility maximization. Elastic services with concave utilities are first considered and the exact expression of optimal bandwidth allocation for each peer is deduced. In order to obtain an optimum in distributed P2P networks, we develop a gradient-based bandwidth allocation scheme and illustrate the performance with numerical examples. Then we investigate bandwidth allocation for inelastic services with sigmoidal utilities, which is a nonconvex optimization problem. In order to solve it, we analyze provider capacity provisioning for bandwidth allocation of inelastic services and modify the update rule for prices that service customers should pay. Numerical examples are finally given to illustrate that the improved scheme can also efficiently converge to the global optimum.

Keywords: P2P networks, bandwidth allocation, elastic and inelastic services, utility function.

1. Introduction

Over the last several years, peer-to-peer (P2P) networks have received much attention due to the fact that they offer a simple way to exchange video files as well as to provide other network services over the Internet. Different from the traditional client-server architecture, which mainly depends on a small number of powerful servers, P2Ps have many advantages such as high scalability and strong robustness. The main idea behind P2Ps is that each peer not only receives resources from the networks also but can provide resources to the networks. A peer can derive much more benefit from the networks if a larger number of peers exchange their resources. Thus, P2Ps have been applied into many fields in recent years, e.g., distributed storage (Yan *et al.*, 2017), cloud computing (Song *et al.*, 2014), edge computing (Wang *et al.*, 2018) and hierarchical name systems (Lin *et al.*, 2015).

P2Ps can support various network services and have caused a lot of Internet traffic by different protocols (Song *et al.*, 2015; 2017; Zheng *et al.*, 2016), such as BitTorrent for file-sharing and VoIP for video conferencing. Indeed, these services can be mainly divided into two categories

according to the satisfaction of peer when acquiring services (Lee *et al.*, 2005; Hande *et al.*, 2007; Li *et al.*, 2015; 2016b). One corresponds to the traditional data services, such as file download and upload. They are almost tolerant of transmission delay and can make use of even the minimal amounts of bandwidth. These services are known as elastic services and the satisfaction of a peer when requesting an elastic service can be described by a concave utility. The other one is related to delay or rate sensitive multimedia services, such as real-time streaming video service. They always have high requirements on time delay and the service rate for receiving a certain level of the quality of service (QoS). These services are regarded as inelastic in their requirement for bandwidth. The satisfaction of each peer for obtaining an inelastic service can be modeled as nonconcave utility (e.g., sigmoidal).

Resource allocation for elastic and inelastic services in P2P networks has been an emerging area of research. Resource pricing is regarded as an interesting means to realize resource allocation. Eger and Killat (2007b) proposed a pricing mechanism to achieve fair bandwidth allocation of service providers between service requesters. They further studied the weighted fairness among

*Corresponding author

service requesters and presented an extended bandwidth allocation mechanism where service requesters adjust their offered prices and service providers adjust their service rates (Eger and Killat, 2007a). Then Kumar *et al.* (2011) developed a scheme for pricing and resource allocation in P2P networks, which permits users in a firm to share computing resources effectively. Koutsopoulos and Iosifidis (2010) investigated bandwidth allocation in a star topology P2P network where the access links to the backbone networks become the capacity bottleneck. The authors formulated the bandwidth allocation problem with the objective of maximizing total network utility through reasonably allocating the bandwidth of each peer to downloaders and uploaders. Recently, Li and Sun (2016) as well as Li *et al.* (2016a) applied the first-order Lagrangian idea and low-pass filtering method, and proposed a novel price-based bandwidth allocation mechanism. Thus a service customer receives its bandwidth allocation according to its offered price, which realizes the goal of fair resource allocation. Antal and Vinkó (2016) considered max-min fair bandwidth allocation in BitTorrent communities and presented a mathematical programming model for max-min fairness bandwidth allocation in a multi-swarm peer-to-peer content sharing community.

Besides the research results mentioned above, there are also some other resource allocation mechanisms based on reputation methods. They are found useful to encourage cooperation amongst selfish peers and realize efficient resource allocation. Satsiou and Tassioulas (2010) presented a distributed reputation-based mechanism to achieve resource allocation according to which peers earn reputation analogous to their contributions. Gupta *et al.* (2015) described a scheme of reputation-based probabilistic resource allocation for avoiding free riding in unstructured P2P networks. Goswami *et al.* (2017) considered reputation-based resource allocation in P2P networks and analyzed a resource allocation mechanism by using two non-cooperative games: the demand game and the reputation game.

The research results above mainly concern resource allocation for elastic services, such as file-transfer in P2P file-sharing networks. However, how to achieve efficient resource allocation for inelastic services is also very important. Chen *et al.* (2012) considered resource allocation for P2P multiparty conferencing applications where it is a crucial challenge to provide a certain level of the QoS. The authors described the quality of experience of the conferencing peer through a utility function and formulated the utility maximization model for P2P multiparty conferencing applications, which are constrained by peers' uplink capacities. Li *et al.* (2017) also developed the utility maximization problem for P2P inelastic services, and derived that the problem is difficult to resolve by using traditional schemes due

to the nonconvexity of the optimization. In order to overcome the difficulty and obtain the optimum, they applied particle swarm optimization (PSO).

In this work we investigate bandwidth allocation for multiclass services, and develop a utility maximization (social welfare) model, i.e., the total satisfaction of all peers in the networks. Our work is different from the results obtained previously. The objective of our model is to maximize the aggregated utility of peers when they acquire services, elastic and/or inelastic. We firstly derive a scheme of gradient-based bandwidth allocation for only elastic services, and extend it to apply into bandwidth allocation for inelastic services by slightly modifying the paid price rule, even though the bandwidth allocation model is a difficult nonconvex optimization problem. The simulation results validate that the scheme can achieve the global optimum within a reasonable number of iterations.

The rest of this paper is organized as follows. We introduce the bandwidth allocation model for multiclass services in P2P networks in Section 2. In Section 3 we analyze the model for only elastic services by nonlinear programming theory and present a gradient-based bandwidth allocation scheme. Then we apply the bandwidth allocation algorithm for inelastic services by slightly modifying the price rule in Section 4, and give some numerical examples in Section 5. Finally, we conclude this paper in Section 6.

2. Bandwidth allocation model

2.1. Services and utility functions. We all know that there are many different services in the Internet. Each user who requires a service can have a certain of satisfaction when the service is guaranteed by a certain amount of resource. A utility function is usually used to describe the satisfaction of a user when he/she obtains a service. Based on the features of their functions, the services can be classified into two types (Lee *et al.*, 2005; Hande *et al.*, 2007; Li *et al.*, 2015). One is known as the traditional data services, such as file download and upload. These services are not sensitive to transmission delay and can make use of even the minimal amounts of granted bandwidth. These are considered to be elastic and they have concave utility functions. The other type is related to multimedia services. These services are very sensitive to transmission delay and usually demand a certain amount of bandwidth to support required QoS. Examples of these services include real-time streaming video and audio services. They are inelastic and usually have nonconcave utility functions, such as sigmoidal ones. Here all utility functions are increasing and no less than zero in their arguments, i.e., $U^s(y) \geq U^s(0) = 0$.

We adopt the utility functions discussed by Lee *et al.* (2005), Hande *et al.* (2007) and Li *et al.* (2015). Elastic service s has the concave utility function $U^s(y) =$

$w(\log(ay+b)+d)$ with service rate y , and inelastic service r has the sigmoidal utility function

$$U^r(y) = w \left(\frac{1}{1 + e^{-a(y-b)}} + d \right),$$

where a , b , d and w are parameters of service s or r . Generally, w is regarded as the willingness-to-pay of the customer who requests service s or r , a and b are considered the elasticity of service s or the inelasticity of service r . For example, the larger the parameters a and b , the more inelastic service r becomes, which means a greater bandwidth requirement to guarantee its QoS.

2.2. Model description. Consider a P2P network which consists of a set of peers, a set S of elastic services and a set R of inelastic services. Each peer in the network intends to acquire one or several services. At the same time it can also provide one or several services. Therefore, each peer acts as not only a service customer, but also as a service provider. For file-sharing P2P networks, each peer uses its access link to upload/download a file or a fragment of a file to other peers, who acquire the file. Hence the scarce resource of this network is the upload capacity of each peer, resulting in competition among the requesting peers. Thus the network faces an important problem of resource management, that is, how to efficiently allocate the bandwidth of peers' access links among the service requesting ones.

Introduce the set P of peers acting as service providers that grant bandwidth allocation to requesters. Also define the sets C^s and C^r of peers acting as service customers that request elastic services and inelastic services, respectively. Peer $c \in C^s$ obtains a total rate y_c^s granted by its providers $P^s(c)$ when requesting elastic service s , and peer $c \in C^r$ obtains a total rate y_c^r granted by its providers $P^r(c)$ when requesting inelastic service r . For peer p as a service provider, assume $x_{pc}^s \geq 0$ is the service rate offered by service provider p for customer c who requests elastic service s and $x_{pc}^r \geq 0$ is the service rate offered by service provider p for customer c that requests inelastic service r . Then the total rate offered by provider p is subjected to its access link upload capacity C_p .

Thus bandwidth allocation for multiclass services in P2P networks can be modeled as the following optimization problem:

Maximize

$$\sum_{c \in C^s} U_c^s(y_c^s) + \sum_{c \in C^r} U_c^r(y_c^r) \quad (1)$$

subject to

$$\sum_{p \in P^s(c)} x_{pc}^s = y_c^s, \quad \sum_{p \in P^r(c)} x_{pc}^r = y_c^r,$$

$$\sum_{c \in C^s(p)} x_{pc}^s + \sum_{c \in C^r(p)} x_{pc}^r \leq C_p.$$

The bandwidth allocation problem (1) can be considered an P2P SYSTEM problem. Here, the objective of bandwidth allocation is to maximize the aggregated utility of service rates y_c^s and y_c^r over all service customers with constraints of service providers' access link capacities. Notice the equality in the bandwidth allocation model; for each service customer c , the elastic (inelastic) service rate y_c^s (y_c^r) is the sum of the rates x_{pc}^s (x_{pc}^r) that its service providers offer. On the other hand, as described by the inequality in the optimization problem, the service rate of provider p is constrained by its upload capacity, i.e., C_p .

2.3. Model analysis. Now we analyze the bandwidth allocation model (1) for multiclass services in peer-to-peer networks. The Lagrangian of the model (1) is

$$\begin{aligned} L(\mathbf{x}, \mathbf{y}; \lambda, \mu) &= \sum_{c \in C^s} U_c^s(y_c^s) + \sum_{c \in C^r} U_c^r(y_c^r) \\ &+ \sum_{c \in C^s} \lambda_s \left(\sum_{p \in P^s(c)} x_{pc}^s - y_c^s \right) \\ &+ \sum_{c \in C^r} \lambda_r \left(\sum_{p \in P^r(c)} x_{pc}^r - y_c^r \right) \\ &+ \sum_{p \in P} \mu_p \left(C_p - \sum_{c \in C^s(p)} x_{pc}^s - \sum_{c \in C^r(p)} x_{pc}^r \right), \end{aligned} \quad (2)$$

where λ is the price vector with elements λ_s and λ_r , which can be considered the price per unit bandwidth paid by customer c when requesting elastic service s and inelastic service r , respectively; μ is the price vector with element μ_p , which can be regarded as the price per unit bandwidth charged by provider p when granting bandwidth allocation for a service; \mathbf{x} is the service rate matrix with elements x_{pc}^s and x_{pc}^r for elastic and inelastic services, respectively; \mathbf{y} is the rate vector with elements y_c^s and y_c^r .

We can rewrite the Lagrangian (2) as

$$\begin{aligned} L(\mathbf{x}, \mathbf{y}; \lambda, \mu) &= \sum_{c \in C^s} (U_c^s(y_c^s) - \lambda_s y_c^s) \\ &+ \sum_{c \in C^r} (U_c^r(y_c^r) - \lambda_r y_c^r) \\ &+ \sum_{c \in C^s} \sum_{p \in P^s(c)} x_{pc}^s (\lambda_s - \mu_p) \\ &+ \sum_{c \in C^r} \sum_{p \in P^r(c)} x_{pc}^r (\lambda_r - \mu_p) + \sum_{p \in P} \mu_p C_p. \end{aligned} \quad (3)$$

Notice that the first part in (3) is separable in y_c^s and y_c^r , and the second part is separable in x_{pc}^s and x_{pc}^r . Thus

the objective function of the dual problem is

$$\begin{aligned}
 D(\lambda, \mu) &= \max_{\mathbf{x}, \mathbf{y}} L(\mathbf{x}, \mathbf{y}; \lambda, \mu) \\
 &= \sum_{c \in C^s} \mathcal{P}^s(\lambda_s) + \sum_{c \in C^r} \mathcal{P}^r(\lambda_r) \\
 &\quad + \sum_{c \in C^s} \sum_{p \in P^s(c)} \mathcal{R}_{pc}^s(\lambda_s, \mu_p) \\
 &\quad + \sum_{c \in C^r} \sum_{p \in P^r(c)} \mathcal{R}_{pc}^r(\lambda_r, \mu_p) + \sum_{p \in P} \mu_p C_p,
 \end{aligned} \tag{4}$$

where

$$\mathcal{P}^s(\lambda_s) = \max_{y_c^s} U_c^s(y_c^s) - \lambda_s y_c^s, \tag{5}$$

$$\mathcal{P}^r(\lambda_r) = \max_{y_c^r} U_c^r(y_c^r) - \lambda_r y_c^r, \tag{6}$$

$$\mathcal{R}_{pc}^s(\lambda_s, \mu_p) = \max_{x_{pc}^s} x_{pc}^s (\lambda_s - \mu_p), \tag{7}$$

$$\mathcal{R}_{pc}^r(\lambda_r, \mu_p) = \max_{x_{pc}^r} x_{pc}^r (\lambda_r - \mu_p). \tag{8}$$

Then, we can obtain the *dual problem* of the bandwidth allocation model (1):

Minimize

$$D(\lambda, \mu) \tag{9a}$$

subject to

$$\lambda_s \geq 0, \quad \lambda_r \geq 0, \quad \mu_p \geq 0. \tag{9b}$$

The dual problem aims to achieve the objective of minimizing the total price charged by all service providers under the constraints that service customers are guaranteed with certain levels of satisfaction. In order to obtain the optimal price and bandwidth allocation, gradient-based schemes could be derived when only considering elastic services since their utility functions are all concave. However, when discussing multiclass services, the bandwidth allocation model becomes a nonconvex problem, such that these traditional gradient-based schemes are not necessarily efficient to converge to the global optimum. They may produce suboptimal or even infeasible bandwidth allocation for each peer.

3. Bandwidth allocation for elastic services

In this part, elastic services with concave utility functions are firstly analyzed through convex optimization theory, and a gradient-based bandwidth allocation scheme for them is proposed for achieving the optimum of the bandwidth allocation model.

3.1. Model analysis. Based on convex optimization theory (Bertsekas, 2003), we find the utility maximization model (1) for bandwidth allocation in P2P networks is

a convex programming problem, since it has a concave objective and linear constraints. The objective of the model (1) is strictly concave with respect to variable y_c^s ; however, it is concave but not strictly concave with respect to variable x_{pc}^s since each customer can receive bandwidth allocation from multiple providers. Thus we can derive the following theorem.

Theorem 1. *For the utility maximization model (1) for bandwidth allocation in P2P networks, if there are only elastic services, then the the optimal aggregated service rate of each customer c , i.e., y_c^{s*} , exists and is unique. Meanwhile, the optimal bandwidth provision of each customer from its providers, i.e., x_{pc}^{s*} , is not necessarily unique.*

If there are only elastic services in the network, the optimal primal objective value is equal to the optimal dual objective value. Analyzing the prices charged by service providers, we also obtain the following result when the bandwidth allocation model (1) achieves its optimum.

Theorem 2. *When the bandwidth allocation model (1) attains its optimum, the prices charges by those providers offering bandwidth allocation to the same customer are all equal, i.e., for service providers $p, q \in P^s(c)$, $p \neq q$, then $\mu_p = \mu_q$. Furthermore, these charged prices are equivalent to the price λ_s offered by the customer.*

Proof. When the bandwidth allocation model (1) attains its optimum, the Karush–Kuhn–Tucker (KKT) condition holds for optimality of the bandwidth allocation problem. Thus, we can obtain

$$U_c^{s'}(y_c^{s*}) - \lambda_s^* = 0 \quad \text{if } y_c^{s*} > 0, \forall c \in C^s,$$

$$\begin{aligned}
 \lambda_s^* - \mu_p^* &= \lambda_s^* - \mu_q^* = 0 \\
 &\quad \text{if } x_{pc}^{s*} > 0, x_{qc}^{s*} > 0, \forall p, q \in P^s(c).
 \end{aligned}$$

Hence, for the multiple providers that grant bandwidth allocation for a customer requesting an elastic service, e.g., $p, q \in P^s(c)$, the optimal price charged by each provider is

$$\mu_p^* = \mu_q^* = \lambda_s^* = \frac{a^s w_c^s}{a^s y_c^{s*} + b^s}. \tag{10}$$

Then the result is obtained. ■

We can also understand this theorem from the relationship of service or bandwidth request and provision. Here μ_p is regarded as the price per unit bandwidth charged by service provider p . It serves as a load indicator on provider p . If the request load on this provider increases, μ_p will increase to show that its resource is “expensive” to use. Thus, μ_p tries to balance the request load by discouraging customers from

requesting bandwidth from those providers where there may occur high load. At the end, the prices charged by providers are all equivalent to that offered by the serving customer.

3.2. Optimal bandwidth allocation. At the optimum of an optimization problem, by applying the KKT condition, the constraint with service provider p in the bandwidth allocation model (1) is *active* if the price associated with it is nonzero, i.e., $\mu_p > 0$; otherwise, the constraint is *nonactive*, which can be ignored in the following analysis.

Notice that the elastic services have strictly concave utility functions. Thus from (5) we derive the optimal aggregated service rate of customer c who requests elastic service s ,

$$y_c^{s*} = U_c^{s'}(\lambda_s)^{-1} = \frac{w_c^s}{\lambda_s} - \frac{b^s}{a^s}. \quad (11)$$

Substituting (11) into (3), we obtain the following simplified Lagrangian:

$$\begin{aligned} \bar{L}(\mathbf{x}; \lambda, \mu) &= \sum_{c \in C^s} \left(w_c^s \log \frac{a^s w_c^s}{\lambda_s} + d^s \right) \\ &\quad + \lambda_s \sum_{p \in P^s(c)} x_{pc}^s - w_c^s + \frac{b^s}{a^s} \lambda_s \\ &\quad + \sum_{p \in P} \mu_p \left(C_p - \sum_{c \in C^s(p)} x_{pc}^s \right). \end{aligned} \quad (12)$$

Let $\partial \bar{L}(\mathbf{x}; \lambda, \mu) / \partial \lambda_s = 0$; then we obtain the optimal price paid by customer c requesting elastic service s ,

$$\lambda_s^* = \frac{w_c^s}{\sum_{p \in P^s(c)} x_{pc}^s + \frac{b^s}{a^s}}. \quad (13)$$

Inserting (13) into (12), we have

$$\begin{aligned} \tilde{L}(\mathbf{x}; \mu) &= \sum_{c \in C^s} w_c^s \left(\log \left(a^s \sum_{p \in P^s(c)} x_{pc}^s + b^s \right) + d^s \right) \\ &\quad + \sum_{p \in P} \mu_p \left(C_p - \sum_{c \in C^s(p)} x_{pc}^s \right). \end{aligned} \quad (14)$$

From (10) and (13), we obtain

$$\begin{aligned} \sum_{p \in P^s(c)} x_{pc}^s &= \frac{w_c^s}{\mu_p} - \frac{b^s}{a^s} \\ &= \frac{w_c^s}{\mu_q} - \frac{b^s}{a^s}, \quad p, q \in P^s(c). \end{aligned} \quad (15)$$

We can construct a bipartite graph which is composed of the two sets C and P . Then an edge denotes

a service relationship between customer c and provider s . If the bipartite graph is fully connected, then $\mu_p^* = \mu_q^* = \lambda_s^*$, $p, q \in P^s(c)$ holds, whereas if the bipartite graph is not connected, a similar optimization process can be run for every disjoint connected subgraph separately. Thus, inserting (15) into (14), the Lagrangian can be rewritten as

$$\begin{aligned} \tilde{L}(\mathbf{x}; \mu) &= \sum_{c \in C^s} w_c^s \left(\log \left(a^s \sum_{p \in P^s(c)} x_{pc}^s + b^s \right) + d^s \right) \\ &\quad + \sum_{p \in P} \mu_p C_p - \sum_{p \in P} \mu_p \sum_{c \in C^s(p)} x_{pc}^s \\ &= \sum_{c \in C^s} \left(w_c^s \left(\log \frac{a^s w_c^s}{\mu} + d^s \right) - w_c^s + \frac{b^s}{a^s} \mu \right) \\ &\quad + \mu \sum_{p \in P} C_p. \end{aligned}$$

Setting $d\tilde{L}(\mathbf{x}; \mu)/d\mu = 0$, we obtain

$$\mu = \frac{\sum_{c \in C^s} w_c^s}{\sum_{p \in P} C_p + \sum_{c \in C^s} \frac{b^s}{a^s}}, \quad (16)$$

and from (15), we have

$$\begin{aligned} y_c^{s*} &= \sum_{p \in P^s(c)} x_{pc}^{s*} \\ &= \frac{w_c^s}{\sum_{c \in C^s} w_c^s} \left(\sum_{p \in P} C_p + \sum_{c \in C^s} \frac{b^s}{a^s} \right) - \frac{b^s}{a^s}. \end{aligned} \quad (17)$$

Thus, if there are only elastic services in the network, the optimal aggregated service rate y_c^{s*} of customer c depends on the sum of customers' parameters b^s/a^s plus the total capacity of all service providers, and the total willingness-to-pay weighted by the willingness-to-pay of the customer. We also find that the optimal aggregated service rate of each customer is unique, as we have discussed in Theorem 1.

3.3. Bandwidth allocation scheme. In order to attain an optimal bandwidth allocation in a decentralized P2P network, distributed schemes for bandwidth allocation depending only on locally available information should be developed. Using the first-order Lagrangian method and a filtering mechanism, we propose the following bandwidth allocation scheme for elastic services. Details of the presented scheme can eliminate typical oscillation behavior due to nonuniqueness of the optimum. The proposed bandwidth allocation scheme are described as follows.

Each service provider p updates its bandwidth allocation for customer c who requests elastic service s with the following scheme:

$$x_{pc}^s(t+1) = ((1-\alpha)x_{pc}^s(t) + \alpha\tilde{x}_{pc}^s(t) + \alpha\kappa x_{pc}^s(t)(\lambda_s(t) - \mu_p(t)))_{x_{pc}^s(t)}^+, \quad (18)$$

$$\tilde{x}_{pc}^s(t+1) = (1-\alpha)\tilde{x}_{pc}^s(t) + \alpha x_{pc}^s(t), \quad (19)$$

where $\kappa > 0$ is the step size; $\alpha > 0$ is the parameter for low-pass filtering in the algorithm, which is used to remove the oscillation and improve the convergence without changing the optimal solution; $a = (b)_c^+$ means $a = b$ if $c > 0$ and $a = \max\{0, b\}$ if $c = 0$.

Each customer c computes the price $\lambda_s(t)$ paid to its providers according to the following rule:

$$\lambda_s(t) = \arg \max_{y_c^s(t)} U_c^s(y_c^s(t)) - \lambda_s(t)y_c^s(t), \quad (20)$$

$$y_c^s(t) = \sum_{p \in P^s(c)} x_{pc}^s(t). \quad (21)$$

Each provider p updates its charged price $\mu_p(t)$ with the following scheme:

$$\mu_p(t+1) = \left(\mu_p(t) + \nu \frac{z_p(t) - C_p}{C_p} \right)_{\mu_p(t)}^+, \quad (22)$$

$$z_p(t) = \sum_{c \in C^s(p)} x_{pc}^s(t), \quad (23)$$

where $\nu > 0$ is the step size.

In the bandwidth allocation algorithm above, customer c obtains the aggregated service rate $y_c^s(t)$ and computes the price $\lambda_s(t)$ paid for provider p according to (20). Provider p updates its bandwidth allocation $x_{pc}^s(t)$ for customer c according to (18), (19). Meanwhile, provider p observes the load $z_p(t)$ on it, and updates its charged price $\mu_p(t)$ according to (22). Obviously, the rule for bandwidth allocation update (18)–(21) and price update (22)–(23) are both gradient-based schemes, which can be proven to converge to the optimum of the convex optimization problem.

When there are inelastic services besides elastic ones in the network, the utility maximization model (1) for bandwidth allocation in P2P networks becomes an intrinsically difficult problem of nonconvex optimization, which will be discussed in the next section.

4. Bandwidth allocation for inelastic services

In this part, we assume that there are both elastic and inelastic services in the P2P network and analyze the utility maximization model (1) for bandwidth

allocation. In this scenario the bandwidth allocation problem becomes a nonconvex optimization because of the sigmoidal utilities of inelastic services. It is hard to obtain the optimum through traditional methods. The gradient-based scheme proposed for bandwidth allocation of elastic services would not necessarily converge for inelastic services. In the following analysis we will show that, if service providers are provisioned with certain amounts of capacity, we can still obtain the optimal bandwidth allocation for both elastic and inelastic services.

4.1. Model analysis. Consider the sigmoidal utility function $U^r(y^r) = w/(1 + e^{-a(y^r-b)}) + d$ of inelastic service r , and construct a straight line from the origin to be tangent to the sigmoidal function. Denote the y -coordinate of the intersection of the tangent with the sigmoidal function by y^{r0} and the slope of the tangent by λ_r^0 , as shown in Fig. 1. We know that there is an inflection point y_0^r for the sigmoidal utility function that satisfies $d^2U^r(y^r)/dy^{r2} > 0$, for $y^r < y_0^r$, and $d^2U^r(y^r)/dy^{r2} < 0$, for $y^r > y_0^r$; that is, the inflection point y_0^r separates the sigmoidal function into two portions; a convex one at a low rate and a concave one at high rate. It can be observed that $y^{r0} \geq y_0^r$ for the sigmoidal function.

We can obtain from (6) that the optimal price offered by customer c who requests inelastic service r can be denoted as $\lambda_r^* = U_c^{r'}(y_c^r)$, which is a function of the aggregated rate y_c^r . For one multiplier λ_r^* , there are two options for inelastic service with a sigmoidal utility function, i.e., $y^{r1}(\lambda_r^*)$ and $y^{r2}(\lambda_r^*)$ ($\geq y^{r1}(\lambda_r^*)$). Indeed, if $\lambda_r > \lambda_r^0$, then $y^{r1*}(\lambda_r) = y^{r2*}(\lambda_r) = \arg \max U^r(y^r) - \lambda_r y^r = 0$; if $\lambda_r \leq \lambda_r^0$, only the larger service rate, i.e., $y^{r2*}(\lambda_r)$, is the optimal bandwidth allocation to maximize the sigmoidal function since the smaller one, i.e., $y^{r1*}(\lambda_r)$, is in fact to minimize the sigmoidal function. In particular, when $\lambda_r = \lambda_r^0$, we get $y^{r2*}(\lambda_r) = y^{r0}$ and $y^{r1*}(\lambda_r) = 0$, since both y^{r0} and 0 maximize $U^r(y^r) - \lambda_r y^r$.

Accordingly, the optimal price offered by customer

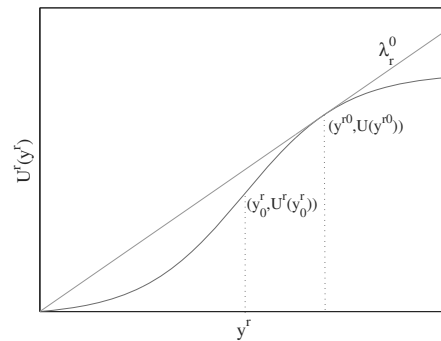


Fig. 1. Sigmoidal utility function and its tangent.

c that requests inelastic service r should be less than λ_r^0 , and the optimal bandwidth allocation for inelastic service can be derived. Furthermore, we can also obtain from the analysis above that $y_c^{r*}(\lambda_r) \geq y_c^{r0}(\lambda_r^0)$ if $\lambda_r \leq \lambda_r^0$ or $\lambda_r \leq \lambda_r^0$ if $y_c^{r*} \geq y_c^{r0}$. Therefore, in order to attain the aim that the price offered by customer c requesting inelastic service r does not exceed the point where the tangent from the origin intersects the sigmoidal function, i.e., λ_r^0 , each inelastic service should be guaranteed with a certain amount of bandwidth allocation which is not smaller than the y -coordinate of the intersection of the tangent with the sigmoidal function, i.e., y_c^{r0} .

Generally speaking, increasing the capacities of service providers can reduce the optimal prices to a certain value so as to realize the optimal bandwidth allocation for both elastic and inelastic services. Next, we will investigate the upload capacities of service providers that grant bandwidth allocation for requesting customers, and derive sufficient conditions for provisioning the global optimum.

4.2. Capacity provisioning. In order to analyze the conditions of capacity provisioning for the global optimum, we regard provider capacity C_p as a function of price μ_p , i.e., $C_p(\mu_p)$. With no loss of generality, we suppose that all providers serve both elastic and inelastic services; otherwise, the providers that do not serve any inelastic service can be ignored in analysis, since increasing the capacity on those providers has no contribution towards improving the QoS of inelastic services.

Theorem 3. Consider the utility maximization model for bandwidth allocation of multiclass services in P2P networks. The optimal bandwidth allocation for both elastic and inelastic services can be derived if there exists a price vector $\mu \geq 0$ with element $\mu_p \geq 0$ which satisfies the inequality $\mu_p < \lambda_r^0, p \in P^r(c)$, and the capacity $C_p(\mu_p)$ for each service provider $p \in P$, where

$$C_p(\mu_p) = \sum_{c \in C^s(p)} x_{pc}^s(\mu_p) + \sum_{c \in C^r(p)} x_{pc}^r(\mu_p). \quad (24)$$

and

$$y_c^s = \sum_{p \in P^s(c)} x_{pc}^s(\mu_p), y_c^r = \sum_{p \in P^r(c)} x_{pc}^r(\mu_p). \quad (25)$$

Here, $x_{pc}^s(\mu_p)$ and $x_{pc}^r(\mu_p)$ are the price-based rate allocation at price μ_p obtained by solving (5)–(8) for services with concave and sigmoidal utility functions, respectively.

Proof. Under the assumption above, the subgradient of the Lagrangian dual function $D(\lambda, \mu)$ with respect to link

price μ_p is provided by

$$\frac{\partial D(\lambda, \mu)}{\partial \mu_p} = C_p - \sum_{c \in C^s(p)} x_{pc}^s - \sum_{c \in C^r(p)} x_{pc}^r.$$

If the link capacities satisfy $C_p = C_p(\mu_p)$, the subgradient at access link price μ_p vanishes, i.e., $\partial D(\lambda, \mu)/\partial \mu_p = 0$. Thus the complementary slackness condition holds and hence μ_p is the dual optimum. Here $x_{pc}^s(\mu_p)$ and $x_{pc}^r(\mu_p)$ constitute indeed the primal optimum. As for the inelastic service with a sigmoidal utility function, $\lambda_r^* = \mu_p, p \in P^r(c)$ and $\lambda_r^0 > \mu_p, p \in P^r(c)$. Thus, $\lambda_r^* < \lambda_r^0$, i.e., the optimal price λ_r^* paid by customer c that requests inelastic service r is smaller than the slope at the critical point where the tangent from the origin intersects the sigmoidal function. Thus, the optimal aggregated bandwidth allocation y_c^r of customer c is larger than the critical service rate y_c^{r0} , the point where the tangent from the origin intersects the sigmoidal utility. The sigmoidal function satisfies $y^{r0} \geq y_0^r$; then the optimum y_c^r lies in the concave part of the sigmoidal curve and the optimum $(x_{pc}^s(\mu_p), x_{pc}^r(\mu_p))$ is global. Since μ_p is the optimal price of the dual problem, $x_{pc}^s(\mu_p)$ and $x_{pc}^r(\mu_p)$ are the bandwidth allocation that service provider p offers for customers requesting elastic and inelastic services, respectively. The proof of this theorem is completed. ■

Notice that the optimal granted bandwidth allocation for customers is nonincreasing with respect to the service provider price μ_p . Thus we can claim that the more inelastic the services (i.e., the smaller λ_r^0), the larger the provider capacity C_p needed for provision of the global optimum.

4.3. Bandwidth allocation scheme. With guarantee of the sufficient condition in Theorem 3, the gradient-based bandwidth allocation scheme (18)–(23) proposed for elastic services in Section 3 could also be applied to bandwidth allocation for multiclass services with a mix of concave and sigmoidal utilities. In order to guarantee the convergence to the global optimum, we modify the price $\lambda_r(t)$ paid by customer c when requesting inelastic service r into the following rule:

$$\lambda_r(t) = \left[\arg \max_{y_c^r(t)} U_c^r(y_c^r(t)) - \lambda_r(t) y_c^r(t) \right]_0^{\lambda_r^0}, \quad (26)$$

$$y_c^r(t) = \sum_{p \in P^r(c)} x_{pc}^r(t). \quad (27)$$

Here, $a = [b]_0^a$ means $a = \min\{c, \max\{0, b\}\}$.

In this modified price rule, the price $\lambda_r(t)$ paid by customer c is not larger than λ_r^0 and the granted bandwidth allocation $y_c^r(t)$ for customer c is not smaller than y_c^{r0} .

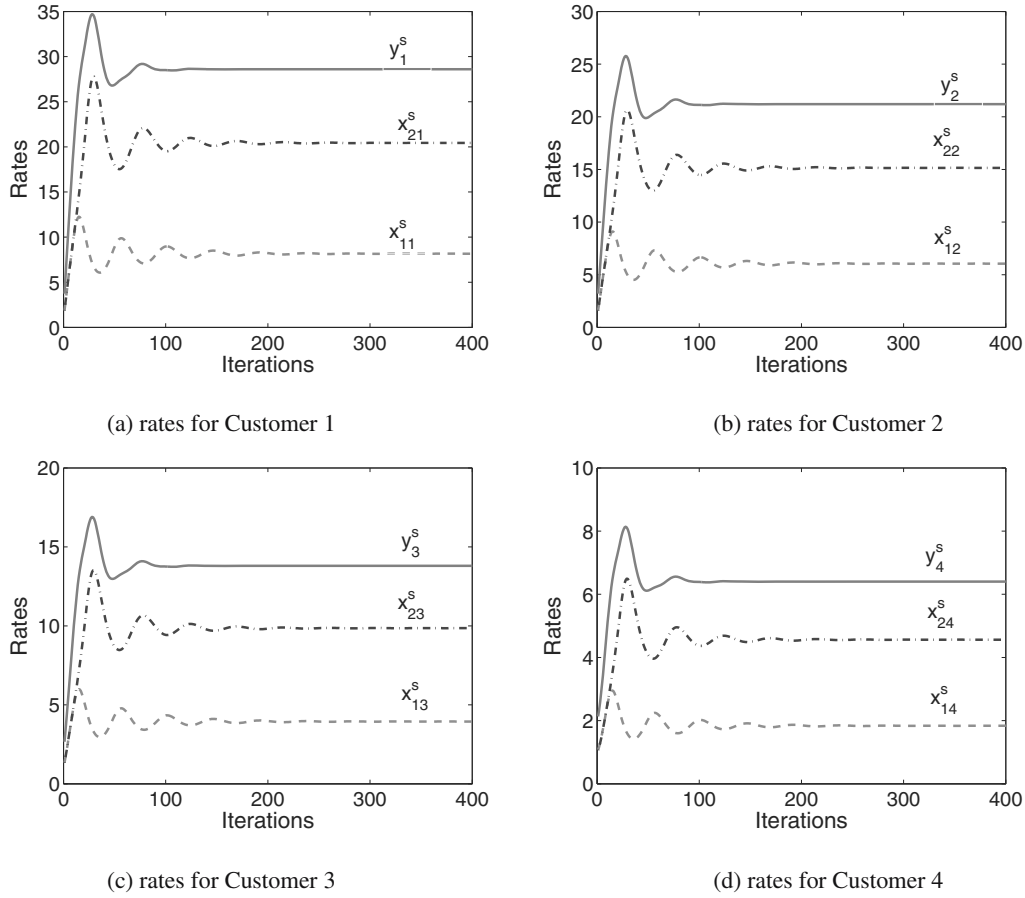


Fig. 2. Optimal bandwidth allocation for elastic services.

Thus the bandwidth allocation scheme converges to the global optimum despite nonconcavity of utility functions, and customers requesting inelastic services attain certain levels of the QoS guarantee.

5. Further discussion

In this section we discuss different forms of utility functions for elastic and inelastic services proposed by Vo et al. (2012). Firstly, we consider the following utility functions for elastic services:

$$U_c^s(y_c^s) = \begin{cases} w_c^s \log(y_c^s + 1) & \text{if } \alpha = 1, \\ w_c^s \frac{(y_c^s + 1)^{1-\alpha} - 1}{1 - \alpha} & \text{if } \alpha > 0 \\ & \text{and } \alpha \neq 1. \end{cases} \quad (28)$$

We find that the utility functions are also concave and the proposed bandwidth allocation scheme (18)–(23) can be applied for elastic services as well. In this case the optimal bandwidth allocation for customer c when

requesting elastic service s is

$$y_c^{s*} = \frac{(w_c^s)^{\frac{1}{\alpha}}}{\sum_{c \in C^s} (w_c^s)^{\frac{1}{\alpha}}} \left(\sum_{p \in P} C_p + |C^s| \right) - 1, \quad (29)$$

which depends on the elastic service parameter α , the number of customers who request this elastic service s , and the total willingness-to-pay of customers weighted by the willingness-to-pay of customer c .

Now we consider the sigmoidal utility functions for inelastic services as follows:

$$U_c^r(y_c^r) = w_c^r \frac{(y_c^r)^\beta}{(y_c^r)^\beta + m}, \quad \forall \beta > 1, m > 0. \quad (30)$$

For this type of sigmoidal utility function $U_c^r(y_c^r)$, there also exists a critical point where the tangent from the origin intersects the sigmoidal function, i.e., λ_r^0 . Bandwidth allocation for customer c when requesting inelastic service r should be guaranteed with a certain value that is not smaller than y_c^{r0} . Then the analysis in Section 4 also holds for this type of utility functions, and the proposed bandwidth allocation scheme with a

modified price update rule (26)–(27) can be applied for inelastic services.

6. Simulation results

In this part we analyze the performance of the proposed bandwidth allocation scheme in a P2P network through some numerical examples. We choose step sizes $\kappa = \nu = 0.2$ and the low-pass filtering parameter $\alpha = 0.2$ for the proposed bandwidth allocation algorithm.

6.1. Elastic services. In this part we assume that the customers are only requesting elastic services and analyze the proposed bandwidth allocation scheme. First we consider a simple P2P network which consists of two service providers and four service customers. The access link capacity of service providers is $C = (C_1, C_2) = (20, 50)$ Mbps. The utility functions of elastic services are given by $U_1^s(y_1^s) = 20 \log(y_1^s + 1)$, $U_2^s(y_2^s) = 15 \log(y_2^s + 1)$, $U_3^s(y_3^s) = 10 \log(y_3^s + 1)$, and $U_4^s(y_4^s) = 5 \log(y_4^s + 1)$.

The simulation results of the bandwidth allocation scheme (18)–(23) are illustrated in Figs. 2 and 3, which show the service rates of each customer granted by providers and the performance of the proposed scheme.

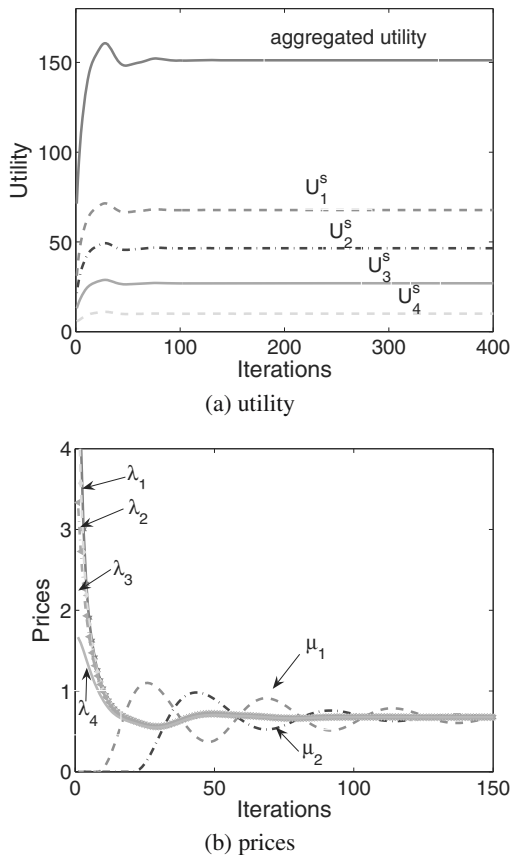


Fig. 3. Performance of the bandwidth allocation algorithm for elastic services.

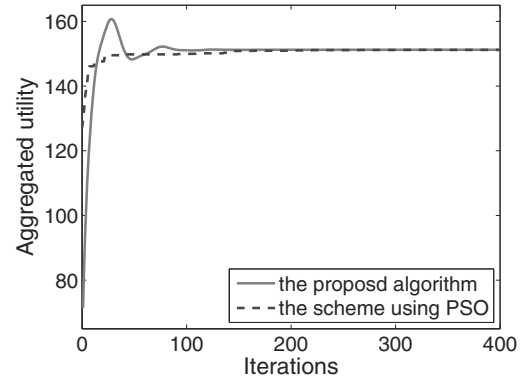


Fig. 4. Comparison of different algorithms for elastic services.

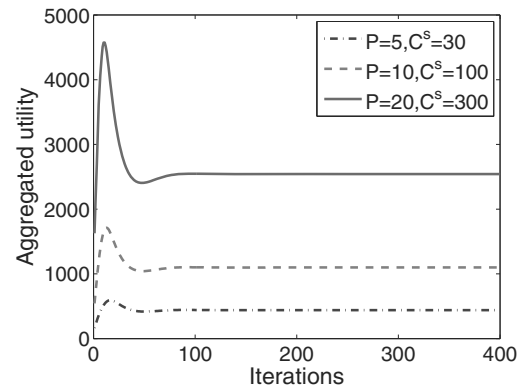


Fig. 5. Performance of the algorithm in different network scenarios for elastic services.

We observe from the simulation results that the bandwidth allocation scheme gradually tends to a steady state where the access link of each provider is approximately 100% utilized. This can be understood from the selfish feature of each peer as customer, acquiring as much resource as possible. The scheme can achieve the optimal bandwidth allocation within reasonable iteration times, providing efficient bandwidth allocation for elastic services. Also, we observe that the prices charged by two providers that offer service to the same customer are both equivalent to the price paid by the customer (i.e., $\lambda_s^* = \mu_p^*$, $p \in P^s(c)$), which is illustrated in Fig. 3(b).

Table 1 lists the optimal bandwidth allocation by using the proposed scheme. Also, it presents the optimal solution solved by the nonlinear programming software LINGO. We can observe that the optimal bandwidth allocation for customers is not unique since the objective of the bandwidth allocation model is not strictly concave with respect to $x = (x_{pc}^s, c \in C, p \in P, s \in S)$. However, we also find that the optimal total bandwidth allocation for each customer is unique, which has been verified in Theorem 1. Meanwhile, in this case the optimal bandwidth allocation for each customer can also be derived from (17), which is equivalent to the optimal values provided in Table 1 (e.g., $y_1^{s*} = 28.6000$ Mbps).

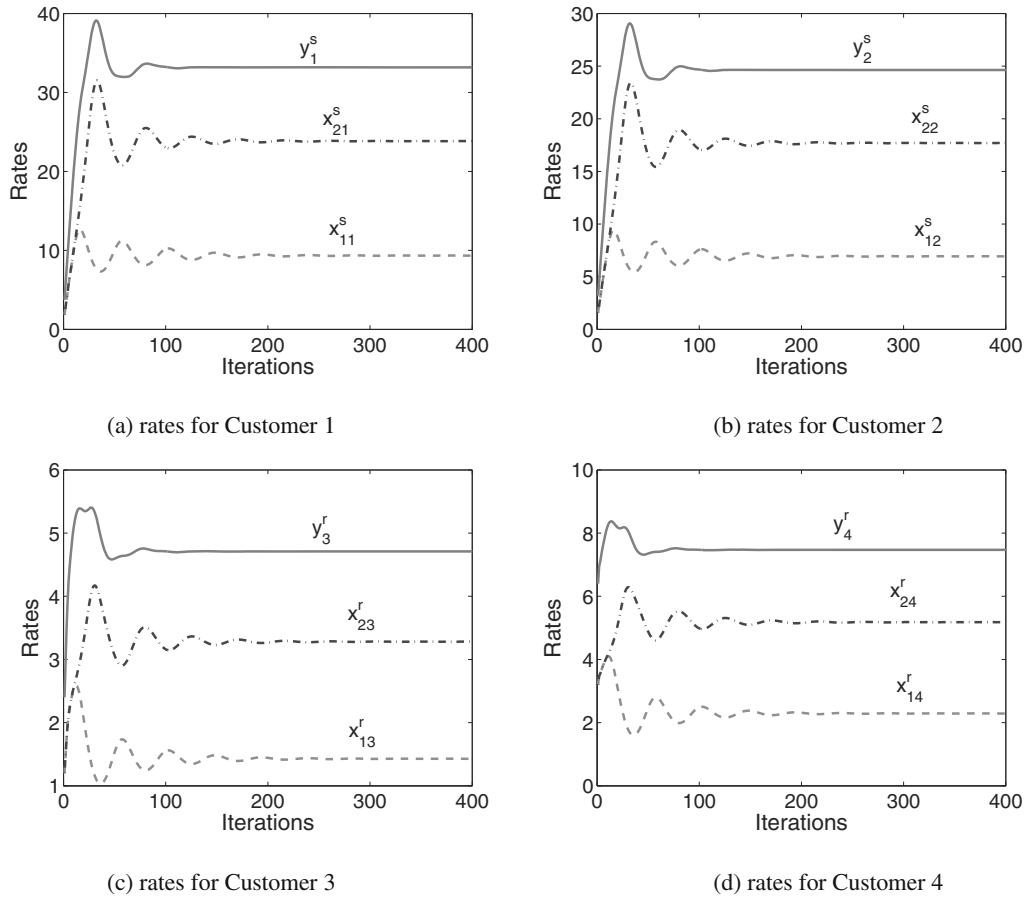


Fig. 6. Optimal bandwidth allocation for elastic and inelastic services.

In addition, we compare the performance of our bandwidth allocation algorithm with the existing schemes such as PSO-based resource allocation (Li et al., 2017), and depict the evolution of the aggregated utility in Fig. 4. In the simulation we choose swarm size 20 for the PSO-based scheme. As we observe in the result, the aggregated utility is gradually driven to the optimal value 151.21 within a reasonable number of iterations for both

the proposed algorithm and the PSO-based scheme.

We also consider the performance of the proposed bandwidth allocation algorithm in different network scenarios. Assume the access link capacity of service providers is 20 Mbps, and the willingness-to-pay of service customers is 10. In Fig. 5, we depict the evolution of the aggregated utility in P2P networks with a different number of peers. We find that the final aggregated utility increases with the number of peers but, in all cases, the optimal value is reached within almost the same number of iterations (e.g., 200). Thus the algorithm can work well in general network scenarios, and the size of the network does not affect its performance.

Table 1. Optimum for the bandwidth allocation model: elastic services.

variable	x_{11}^{s*}	x_{21}^{s*}	x_{12}^{s*}	x_{22}^{s*}
algorithm	8.1650	20.4350	6.0536	15.1464
LINGO	10.8238	17.7762	8.5618	12.6382
variable	x_{13}^{r*}	x_{23}^{r*}	x_{14}^{r*}	x_{24}^{r*}
algorithm	3.9439	9.8561	1.8376	4.5624
LINGO	0.0727	13.7273	0.5417	5.8583
variable	y_1^{s*}	y_2^{s*}	y_3^{s*}	y_4^{s*}
algorithm	28.6000	21.2000	13.8000	6.4000
LINGO	28.6000	21.2000	13.8000	6.4000

6.2. Elastic and inelastic services. Now, in this section we consider bandwidth allocation for multiclass services, that is, the services requested by the customer are not only elastic but also inelastic, and investigate the performance of the proposed bandwidth allocation scheme. Also we first consider a simple P2P network which consists of two service providers with access link capacity $C = (C_1, C_2) = (20, 50)$ Mbps and four service customers. The first two customers are requesting elastic

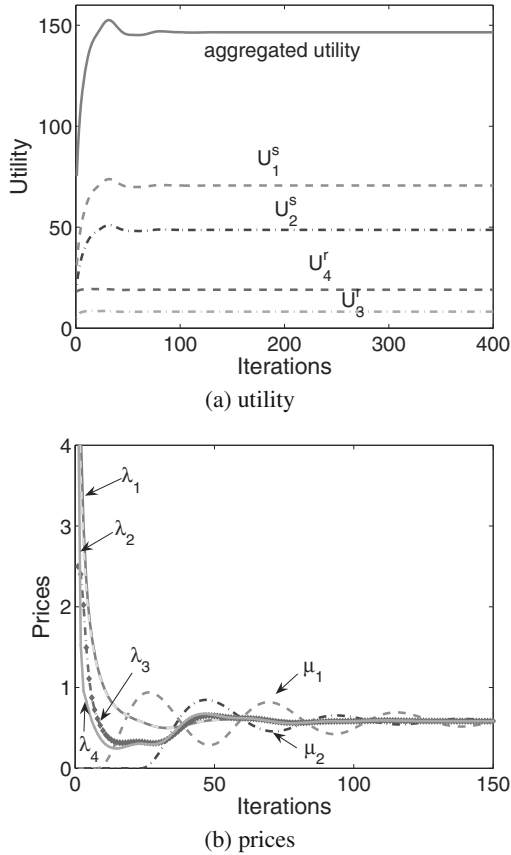


Fig. 7. Performance of the bandwidth allocation algorithm for elastic and inelastic services.

services, and the utility functions of elastic services are given by $U_1^s(y_1^s) = 20 \log(y_1^s + 1)$ and $U_2^s(y_2^s) = 15 \log(y_2^s + 1)$. The other customers are requesting inelastic services and the utility functions of inelastic services are

$$U_3^r(y_3^r) = 10 \left(\frac{1}{1 + e^{-(y_3^r - 2)}} - \frac{1}{1 + e^2} \right)$$

$$U_4^r(y_4^r) = 20 \left(\frac{1}{1 + e^{-(y_4^r - 4)}} - \frac{1}{1 + e^4} \right).$$

We can obtain that, for these two inelastic services, $y_3^{r0} = 2.9166$ Mbps, $\lambda_3^0 = 2.045$, and $y_4^{r0} = 5.5398$ Mbps, $\lambda_4^0 = 2.9078$.

We obtain the simulation results of the proposed bandwidth allocation scheme (18)–(23) with the modified price rule (26), (27) for inelastic services and illustrate them in Figs. 6 and 7. Obviously, we observe that the scheme can converge to the optimal bandwidth allocation for multiclass services within a reasonable number of iterations, i.e., $x^* = (9.3432, 23.8396, 6.9383, 17.6987, 1.4280, 3.2817, 2.2903, 5.1802)$ Mbps. Meanwhile, we also find from Fig. 7(b) that the optimal price paid by customers is $\lambda_r^* = 0.5851$, which is smaller than the

critical prices $\lambda_3^0 = 2.045$ and $\lambda_4^0 = 2.9078$ of the two inelastic services. The prices charged by the two providers that offer services are both equal to the price paid by each customer. In this case, the condition derived from Theorem 3 is satisfied. Thus the optimum is global.

We also compare the performance of our bandwidth allocation algorithm with that of a PSO-based resource allocation scheme (Li *et al.*, 2017) for inelastic services, and depict the evolution of the aggregated utility in Fig. 8. Similarly to the case for only elastic services, the aggregated utility is gradually driven to the optimal value of 146.50 within a reasonable number of iterations for both the proposed algorithm and the PSO-based scheme. We also investigate the robustness of the proposed algorithm in different network scenarios.

Assume that the willingness-to-pay of customers for elastic services is 20 and for inelastic services it is 10. The key parameters are $a = 1, b = 2$ for inelastic services. In Fig. 9, we depict the evolution of the aggregated utility for both elastic and inelastic services in P2P networks with a different number of peers. As we can observe from the simulation, the proposed algorithm works well to achieve the optimal bandwidth allocation for multiclass services in P2P networks with a different number of peers.

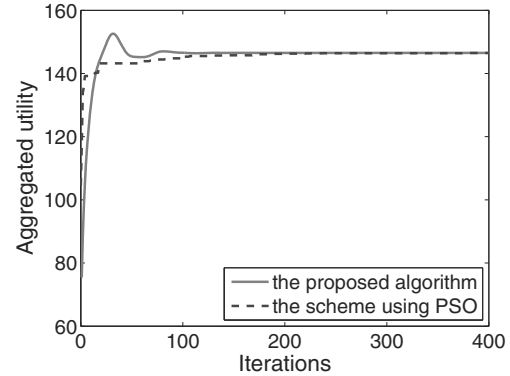


Fig. 8. Comparison of different algorithms for elastic and inelastic services.

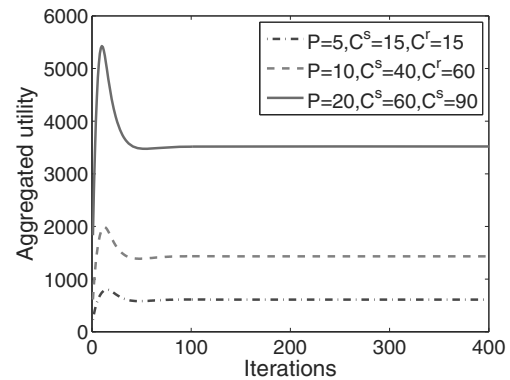


Fig. 9. Performance of the algorithm in different network scenarios for elastic and inelastic services.

7. Conclusions

In this paper we considered the optimal bandwidth allocation for both elastic and inelastic services in P2P networks, and formulated the utility maximization model for peers who request these services. First we considered only elastic services with concave utilities, and obtain the expression of the optimal bandwidth allocation of each peer. To attain an optimum in decentralized networks, we developed a gradient-based bandwidth allocation scheme. However, the scheme may not work well for bandwidth allocation for inelastic services due to the non-concavity of utilities. In order to overcome this, we discussed the capacity provisioning for bandwidth allocation of inelastic services and modified the update rule for prices that customers should pay. Some numerical examples were finally given to verify the performance of the bandwidth allocation scheme for both elastic and inelastic services.

Acknowledgment

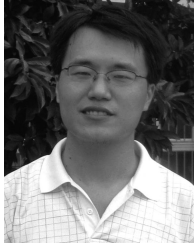
This work was supported in part by the National Natural Science Foundation of China (71671159, 71301139), the Humanity and Social Science Foundation of the Ministry of Education of China (16YJC630106), the Natural Science Foundation of Hebei Province (G2018203302, G2016203236), the project funded by the Hebei Education Department (BJ2017029, BJ2016063) and the project funded by the Hebei Talents Program (A2017002108).

References

- Antal, E. and Vinkó, T. (2016). Modeling max-min fair bandwidth allocation in bittorrent communities, *Computational Optimization and Applications* **66**(2): 383–400.
- Bertsekas, D. (2003). *Nonlinear Programming*, Athena Scientific, Belmont, MA.
- Chen, M., Ponec, M., Sengupta, S., Li, J. and Chou, P. (2012). Utility maximization in peer-to-peer systems with applications to video conferencing, *IEEE/ACM Transactions on Networking* **20**(6): 1681–1694.
- Eger, K. and Killat, U. (2007a). Fair resource allocation in peer-to-peer networks (extended version), *Computer Communications* **30**(16): 3046–3054.
- Eger, K. and Killat, U. (2007b). Resource pricing in peer-to-peer networks, *IEEE Communications Letters* **11**(1): 82–84.
- Goswami, A., Gupta, R. and Parashari, G. (2017). Reputation-based resource allocation in P2P systems: A game theoretic perspective, *IEEE Communications Letters* **21**(6): 1273–1276.
- Gupta, R., Singha, N. and Singh, Y. (2015). Reputation based probabilistic resource allocation for avoiding free riding and formation of common interest groups in unstructured P2P networks, *Peer-to-Peer Networking and Applications* **9**(6): 1101–1113.
- Hande, P., Zhang, S. and Chiang, M. (2007). Distributed rate allocation for inelastic flows, *IEEE/ACM Transactions on Networking* **15**(6): 1240–1253.
- Koutsopoulos, I. and Iosifidis, G. (2010). A framework for distributed bandwidth allocation in peer-to-peer networks, *Performance Evaluation* **67**(4): 285–298.
- Kumar, C., Altinkemer, K. and De, P. (2011). A mechanism for pricing and resource allocation in peer-to-peer networks, *Electronic Commerce Research and Applications* **10**(1): 26–37.
- Lee, J., Mazumdar, R. and Shroff, N. (2005). Non-convex optimization and rate control for multi-class services in the internet, *IEEE/ACM Transactions on Networking* **13**(4): 827–840.
- Li, S., Jiao, L., Zhang, Y., Wang, Y. and Sun, W. (2017). A scheme of resource allocation for heterogeneous services in peer-to-peer networks using particle swarm optimization, *IAENG International Journal of Computer Science* **44**(4): 482–488.
- Li, S. and Sun, W. (2016). A mechanism for resource pricing and fairness in peer-to-peer networks, *Electronic Commerce Research* **16**(4): 425–451.
- Li, S., Sun, W., E, C.-G. and Shi, L. (2016a). A scheme of resource allocation and stability for peer-to-peer file-sharing networks, *International Journal of Applied Mathematics & Computer Science* **26**(3): 707–719, DOI: 10.1515/amcs-2016-0049.
- Li, S., Sun, W. and Hua, C. (2016b). Optimal resource allocation for heterogeneous traffic in multipath networks, *International Journal of Communication Systems* **29**(1): 84–98.
- Li, S., Sun, W. and Tian, N. (2015). Resource allocation for multi-class services in multipath networks, *Performance Evaluation* **92**(1): 1–23.
- Lin, F., Zhou, X., Huang, D. and Yuan, J. (2015). Hierarchical name system based on hybrid p2p for multimedia networks, *Telecommunication Systems* **59**(3): 393–400.
- Satsiou, A. and Tassiulas, L. (2010). Reputation-based resource allocation in p2p systems of rational users, *IEEE Transactions on Parallel and Distributed Systems* **21**(4): 466–479.
- Song, F., Huang, D., Zhou, H., Zhang, H. and You, I. (2014). An optimization-based scheme for efficient virtual machine placement, *International Journal of Parallel Programming* **42**(5): 853–872.
- Song, F., Li, R. and Zhou, H. (2015). Feasibility and issues for establishing network-based carpooling scheme, *Pervasive and Mobile Computing* **24**(1): 4–15.
- Song, F., Zhou, Y., Kong, K., Zheng, Q., You, I. and Zhang, H. (2017). Smart collaborative connection management for identifier-based network, *IEEE Access* **5**: 7936–7949.
- Vo, P., Lee, S. and Hong, C. (2012). The random access num with multiclass traffic, *EURASIP Journal on Wireless Communications and Networking* **242**: 1–12.
- Wang, K., Yin, H., Quan, W. and Min, G. (2018). Enabling collaborative edge computing for software defined vehicular networks, *IEEE Network* **32**(5): 112–117.

Yan, H., Gao, D., Su, W., Foh, C., Zhang, H. and Vasilakos, A. (2017). Caching strategy based on hierarchical cluster for named data networking, *IEEE Access* **5**: 8433–8443.

Zheng, Y., Lin, F., Yang, Y. and Gan, T. (2016). Adaptive resource scheduling mechanism in P2P file sharing system, *Peer-to-Peer Networking and Applications* **9**(6): 1089–1100.



Shiyong Li received his PhD degree from Beijing Jiaotong University, China, in 2011. Currently he is an associate professor in the School of Economics and Management at Yanshan University. He is a (co)author of more than 50 papers in mathematics, engineering, and management journals. He has been a principal investigator/co-investigator in several research projects supported by the National Natural Science Foundation of China, the National Education Committee Foundation of China, the China Postdoctoral Science Foundation, and others. His research interests include resource allocation in networks, electronic commerce, and economics of queues.



Yue Zhang received her BSc degree from Yanshan University, Qinhuangdao, China, in 2017. She is currently working toward her MSc degree in the School of Economics and Management at Yanshan University. Her research interests include resource allocation in cloud computing and electronic commerce.



Yan Wang received her BSc degree from Yanshan University, Qinhuangdao, China, in 2017. She is currently working toward her MSc degree in the School of Economics and Management at Yanshan University. Her research interests include resource allocation in networks and economics of queues.



Wei Sun received her PhD degree from Yanshan University, Qinhuangdao, China, in 2010. Currently she is a full professor in the School of Economics and Management at Yanshan University. She has published more than 40 papers in leading international journals in the areas of operations research and applied mathematics. She has been involved in several projects supported by the National Natural Science Foundation of China, the National Education Committee Foundation of China, and others. Her research interests include economics of queues and queueing systems with vacations.

Received: 14 March 2018

Revised: 3 September 2018

Accepted: 22 October 2018