

Optimal Allocation for Display Advertising

Huaxia Rui,* De Liu[†] and Andrew Whinston[‡]

April 1, 2012

Working Paper: Please do not circulate without permission

Abstract

The spread of real-time targeting technologies in mobile and Internet display advertising creates a new challenge: how to efficiently allocate countless categories of impressions in real time. We argue that existing options for allocating advertising impressions, namely via long-term guaranteed contracts and real-time spot markets, are inefficient. We propose a new “contingent contract” based approach, in which a provider can offer each ad buyer variable quantities contingent on the realized number of impressions. The key advantages of this approach include: (1) it accommodates a wide range of risk preferences, (2) one can compute global optimal contingent contracts ahead of time to achieve better allocation efficiency, and (3) it optimizes the allocation of multiple related impression categories jointly. We first solve the optimal contingent contracts for a single type of impressions and then for multiple categories. We show that a multi-type problem can be converted to many single-type ones, thus allowing efficient computation of optimal contingent contracts. We also present the algorithms for our approach and ways of implementing the optimal allocation in an online environment.

Keywords: display advertising, contingent contracts

1 Introduction

Display advertising is a form of online advertising embedded into a web page, typically include image, text, video, and/or interactive elements designed to convey a marketing message

*University of Texas at Austin, McCombs School of Business, Austin, T e-mail: huaxia@utexas.edu

[†]University of Kentucky, Department of Marketing and Supply Chain, Gatton College of Business and Economics, Lexington, KY 40506, e-mail: de.liu@uky.edu

[‡]University of Texas at Austin, McCombs School of Business, Austin, T e-mail: abw@uts.cc.utexas.edu

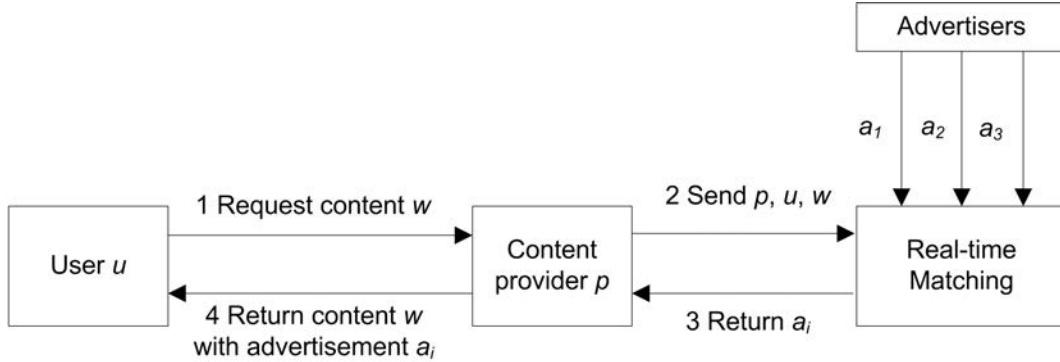


Figure 1: Real-time Targeting in Display Advertising

A user u requests content c from provider p . Provider p forwards information about the user u , content c , and provider p to a matching mechanism. The matching mechanism returns a chosen ad a , which is then served to the user u together with the requested content c .

or cause the user to take an action. Display advertising has seen a dramatic growth in the last 5 years, fueled in part by the popularity of social media, online video, and mobile Internet. A recent forecast suggests that Internet and mobile display advertising will become the largest segment of online advertising by 2015 (EMarketer, 2011).

The market for display advertisement has gone through dramatic changes in the last several years. The most fundamental change is that the industry is moving away from “static” display ads toward dynamic, hyper-targeted display ads based on real-time information (see Figure 1 for a process sketch). Because of availability of real-time information, advertisers can now target users in much smaller granularity, resulting in an explosive number of targetable categories. Facebook, for example, allows display advertisers to target ads based on user’s location, age, gender, relationship status, social network connection types, and interests. McAfee and Papineni (2010) estimated that there are over a trillion distinct categories in display advertising. Contributing to this trend of hyper-targeting is the increasing availability of fine location information: most smart phones are equipped with GPS; indoor mapping technology allows one to pinpoint the location of a mobile user within a building, such as a shopping mall. These new location capabilities enable innovative location-based services as well as location-based display advertisements. Also contributing to this trend of hyper-targeting is the capacity and desire to target users in real time. Cloud computing makes it possible to predict user interests moment by moment and to sniff time-sensitive advertising opportunities based on accumulated data on user behavioral patterns and location information.

Hyper-targeting posts new challenges for the design of display advertising markets. Traditionally, display advertising is sold in bulk via long-term guaranteed contracts that are

negotiated ahead of time (e.g., a year in advance). A typical contract specifies the amount of guaranteed impressions and targeting requirements (e.g., front page, male users, and New York region). If the realized amount of impressions falls below the guaranteed amount, the advertising provider will pay a penalty. This manual, infrequent approach clearly is not scalable to the explosive number of targetable categories we have today; nor does it allow advertisers to take advantage of real-time opportunities. More important, this approach offers undifferentiated guarantees to heterogeneous buyers who may have more or less tolerance for risks. At time of shortage or surplus, it provides no guidance on how resources should be efficiently rationed.

A newer approach relies on real-time spot markets for allocating ad impressions. Facebook, for example, runs a real-time auction for display advertising on its website. The industry is also promoting real-time ad exchanges that are open to multiple buyers and sellers. Yahoo!'s RightMedia Exchange, Google's DoubleClick Exchange, and Microsoft's AdECN are examples of such ad exchanges (see McAfee 2011 for a discussion). In these spot markets, ad buyers can bid on impressions in real time; an automated auction or exchange house allocates the impressions based on the "bids" submitted by the market participants. In theory, spot markets can efficiently allocate display advertising, as long as the market prices vary continuously to balance supply and demand. In practice, however, this approach may incur several inefficiencies.

- First, advertising opportunities are nonstorable. This means that the spot market will operate inefficiently before it finds an efficient market clearing price for the current market condition.
- Second, as argued by Wilson (1989), spot markets are expensive to operate. As the select of advertisement must be completed in a fraction of a second, the spot market and its participants must rely on high-performance computing and networking facility to function. It is also expensive for market participants to continuously monitor prices, adapt their bidding strategies, and execute those strategies in timely fashion.
- A third source of inefficiency may be caused by lack of coordination across correlated auctions. Most spot markets hold a separate auction for each category of advertising impressions. In many cases, however, the demand for different categories are correlated. Consider the following example: both an insurance retailer and local restaurant wish to advertise to a smart phone user who is walking from A to B. The restaurant is willing to pay \$4 only if the user is at A. The insurance retailer is willing to pay \$5 regardless of the location of the user. The insurance retailer would win a bid at A (for a price of 4) and the restaurant would end up nothing – an inefficient allocation. Such

allocative inefficiencies can also happen when bidders are imperfectly informed about others bidders preferences for correlated categories (as defined by location, time, type of users, or other conceivable dimensions).

Due to the limitation of existing approaches, we propose an alternative contingent-contract based approach for allocating display advertising resources. A *contingent contract* is a type of forward contract that specifies the number of impressions allocated contingent on the total realized supply. For example, a contingent contract may allocate 100% of the portion of total realized supply below 10,000 and 25% of the portion above. Another contingent contract may allocate 75% of the portion above 10,000. In this way, we can use contingent contracts to prioritize the allocation of impressions across ad buyers.

A contingent contract approach gathers buyer preferences in advance and determine an optimal contingent contract to offer to each buyer and the corresponding prices. Because the contingent contract approach distinguishes impressions by priorities, an opportunistic buyer will be given low-priority impressions at a discount price. This approach improves allocative efficiency compared to the guaranteed contract approach because now each buyer gets a customized bundle of impressions at different priority levels that closely matches his preference.

The contingent contract approach is less expensive. For ad buyers, they only need to make periodical choices over holistic contingent contracts rather than to engage in continuous monitoring of spot market prices and bid management. For market operators, the calculation of optimal contingent contracts is done offline, which is less expensive than real-time computation.

Perhaps the most important advantage of the contingency contract approach comes from the allocative efficiency gains. The contingent contract approach learns buyer preferences across different categories of impressions in advance. A global optimization can avoid the pitfalls of the spot market approach as illustrated by the insurance and restaurant example. Moreover, it also avoids the costly price discovery process in spot markets since the optimal allocation is determined offline in advance.

A significant feature of our approach is that it explicitly accommodates preferences across multiple categories. Aside from the efficiency gains from jointly allocating multiple categories, this feature also brings about a few additional benefits. One benefit is that we now allow ad buyers to explicitly express the substitutability (or lack of) across multiple categories of impressions, which is not part of the guaranteed contracts or spot markets. Another benefit is that, as a byproduct of solving the optimal allocation problem, we produce reports on which categories should be allocated together and how they should be priced in relation to one another. Such reports provide useful feedback to buyers and help them discover values

of existing or new categories over time.

The contingent contract approach is based on the suppositions that buyers have heterogeneous risk tolerance and substitutability across different categories, and that there is significant efficiency gains by considering these heterogeneities in a joint allocation problem. When these suppositions hold true and preferences are sufficiently persistent, then the contingent contract approach may be a more efficient form of market organization than spot markets and guaranteed contracts.

The notion of contingent contract is not new. Wilson (1989) proposes a priority-based contract for efficiently rationing electricity. Contingent contracts have been used in insurance policies, financial security and derivatives. Economics and finance literatures have extensive discussions on contingent contracts and their designs (Raviv, 1979; Allen and Gale, 1989). But they mostly focus on single category of resources, cash. This paper departs from previous studies by studying optimal contingent contract design for heterogeneous resources. An important contribution of this paper is the development of theory and algorithm for contingent contract design when the uncertain resources are heterogeneous.

We expect the contingent contract approach to be useful in two scenarios. First, it may be used by an advertising provider – a publisher or a publisher network – to sell its advertising inventories directly to ad buyers. Second, it may be used by intermediaries who buy advertising inventories in the spot markets on behalf of their advertising buyer clients. In the former scenario, contingent contract is a substitute for a real-time spot market which may be too expensive or inefficient. In the latter scenario, contingent contract is a complement of the spot market as it allows smaller, less sophisticated buyers to buy impressions through contingent contracts instead of continuous bidding.

This paper focuses on formulating and solving the optimal contingent contract problem for displaying advertising. We structure the rest of our paper as follows: next we review the related literature, followed by a description of our research problem. We solve the optimal contingent contracts for a single category in Section 4. With mild assumption on valuation functions for impressions, we can completely solve the single-category case. In section 5, we extend our analysis to the multiple-category case. We show that, for any given total supply, the optimal allocation problem can be solved as a series of “constrained” problems each of which is bounded by an “indicator matrix.” Using the special structure of the indicator matrices, we can decompose the constrained problem into several subproblems, each of which can be converted into a standard single-category problem. In this way, we turn the original optimal allocation problem into a problem of matrix search. Based on our theoretical results, we describe and implement a novel algorithm for solving the multi-category case; we also demonstrate its effectiveness with numerical experiments. In section 7, we discuss problems

associated with implementing the optimal contingent contracts in a real time environment. Section 8 concludes the paper.

2 Related Literature

Display Advertising. Although display advertising is not new, the real-time allocation of display advertising slots has emerged only in the last few years. Thus, the literature on display advertising is still very nascent. One issue introduced by the real-time spot market is how ad sellers should allocate the resources between guaranteed contracts and spot markets. McAfee and his colleagues (McAfee and Papineni, 2010; Ghosh et al., 2009; McAfee, 2011) at Yahoo! noted potential “cherry-picking” concerns—that ad providers would be tempted to sell high quality impressions in the spot market and leave guaranteed contract buyers with low quality impressions. To address such a concern, they proposed “maximal representativeness”—the idea that if spot market prices are indicators of quality, insisting on a representative sample of impressions at all price levels will prevent ad sellers from selling low-quality impressions to guaranteed contract buyers. Chen (2009) examined the ad seller’s dynamic trade-off between extracting revenue from spot markets and avoiding paying penalty on guaranteed contracts. The researchers in the aforementioned studies did not attempt to optimize risk allocation; rather they took the current selling mechanisms (i.e., guaranteed contracts and spot markets) as given. In contrast, the approach we take in this paper can accommodate a wide range of risk preferences. We also consider the multi-category allocation to further enhance the allocation efficiency.

Sponsored Search Auctions. Our research is closely related to research on sponsored search. Previous research in sponsored search has examined the design of sponsored search auctions for ranking advertisers and determining payments (e.g., Edelman et al., 2007; Varian, 2007; Liu and Chen, 2006; Liu et al., 2010; Athey and Ellison, 2011), and bidding behavior in such auctions (e.g., Ghose and Yang, 2009; Animesh et al., 2009). Although researchers in these studies looked at optimal allocation, most of them assumed unit (or linear) demand and thus did not consider risk preferences. One exception is Chen et al. (2009), which modeled ad buyers’ demand for impressions using a concave function (thus buyers demonstrate risk aversion) and solved the optimal allocation problem under an incomplete information setting. In this paper, we extend Chen et al. (2009) in two significant ways: First, we allow the ad buyers to have different risk preferences. Second, we also consider optimal allocation of multiple categories of impressions.

Ad Allocation with Budget Constraints. Our research is also related to studies dealing with ad allocation problems with budget constraints (Mehta et al., 2007; Aggarwal

et al., 2010). Although this literature does not model risk preference, budget constraints may be used as an (imperfect) instrument for risk management. For example, Mehta et al. (2007) considered an online matching problem in which a stream of heterogeneous impressions were assigned to n advertisers with the goal of maximizing revenue while respecting advertisers’ bids and budgets. Because bidders with budgets can be viewed as a special case of risk-averse utility functions, insights obtained from this paper may also have implications for allocation with budgets.

3 Model Setup

We follow the “optimize-and-dispatch” approach (Parkes and Sandholm, 2005) for the allocation of heterogeneous impressions among market agents. First, the seller solicits information from buyers regarding their willingness to pay for each category of impression. Then, the seller performs offline global optimization, taking into account of all agents’ preferences. This global optimization module generates a contingent contract for each agent and a corresponding price. Then a dispatch module dynamically allocates the available impressions to each agent according to the allocation rule and the supply of impressions.

The main purpose of this paper is to solve the optimization allocation problem. We achieve this in two steps: first by looking at the problem in a single-category setting. Building on that, we develop a theory and algorithm to optimally allocate multi-category impressions.

Let $\mathcal{N} = \{1, 2, \dots, N\}$ be a set of N agents, including one seller, s , and $N - 1$ buyers. The seller can be interpreted as a large publisher, a publisher network (e.g. Google’s DoubleClick for Publishers network), or an agent who buys impressions on behalf of a set of advertisers. The buyers are advertisers or advertising networks. In each case the seller faces the problem of allocating impressions among multiple buyers.

Let $\mathcal{M} \equiv \{1, 2, \dots, M\}$ denote the set of M distinct impression categories and w denote a realized supply vector with element w_m being the realized supply of category m within a given planning horizon. Denote $\xi_i = \begin{pmatrix} \xi_{i1} & \xi_{i2} & \dots & \xi_{im} \end{pmatrix}$ as a *contingent contract* for agent i with element ξ_{im} being a functional mapping from the supply vector to the quantity of category m allocated to agent i .

Let $x_{im} = \xi_{im}(w)$ denote the quantity of category m allocated to agent i when the supply vector is w . We require that

$$\sum_{i \in \mathcal{N}} x_{im} \leq w_m, \forall m \in \mathcal{M} \tag{1}$$

$$x_{im} \geq 0, \forall i \in \mathcal{N}, m \in \mathcal{M} \tag{2}$$

Denote $x_i = (x_{i1} \ x_{i2} \ \dots \ x_{iM})$ and $x = (x_1 \ x_2 \ \dots \ x_N)$ respectively. We assume that the von Neumann-Morgenstern utility function for each agent takes the following form.

$$U_i(x_i, p_i) = u_i(x_i) - p_i = Q_i \left(\sum_{m=1}^M \alpha_{im} x_{im} \right) - p_i, \quad \forall i \in \mathcal{N} \quad (3)$$

where Q_i is the agent i 's *valuation function*, p_i is his payment, and $\alpha_{im} \geq 0$ is his valuation coefficient for category m . We assume Q_i is continuously differentiable, strictly increasing, and strictly concave.¹ We also assume $Q_i(0) \geq 0$ for all $i \in \mathcal{N}$.² $Q_i(0)$ is interpreted as agent i 's *reservation utility* – that is, the utility i can get if he does not participate in this market.

Before continuing, we have a few remarks on this specification of the utility function.

- The concavity of the valuation function captures an agent's risk tolerance: As the concavity of Q_i increases, agent i becomes less tolerant to uncertainty in the number of impressions allocated.
- If coefficient $\alpha_{im} > 0$, then category m is *desirable* for agent i .
- Each agent views all desirable categories as substitutes. The ratio between coefficients corresponds to the marginal rate of substitution. This may seem a limitation, as one can argue that an agent may prefer an assortment of categories. The preference for assortments can be easily captured by an extension of this specification: as illustrated by (4), instead of having one set of substitutable categories, an agent can have k sets of them: categories within each set are substitutes, but categories across sets are not.

$$U_i(x_i, p_i) = Q_{i1} \left(\sum_{m \in \mathcal{M}_1} \alpha_{im} x_{im} \right) + Q_{i2} \left(\sum_{m \in \mathcal{M}_2} \alpha_{im} x_{im} \right) \dots + Q_{ik} \left(\sum_{m \in \mathcal{M}_k} \alpha_{im} x_{im} \right) - p_i, \quad \forall i \in \mathcal{N} \quad (4)$$

The more substitution sets an agent has, the stronger the agent's preference for assortment. Hence the utility function specification (4) captures both substitutability and assortment preferences. Technically, we may treat an agent with k substitution sets as k separate agents and the remainder of our analysis will apply.³ For expositional simplicity, we maintain (3) as the utility function.

¹Our model can be extended to the situation where Q_i is linear for some agents at the cost of messier notations. We do not pursue this generality for readers' convenience. One may also argue that linear forms of Q_i can be approximated by strictly concave Q_i with very small concavity.

²The seller's valuation may be derived from self-use, selling the impressions in a spot market, or simply leaving the advertising slots empty.

³The optimal allocation problem is scalable in the number of agents. So this extension neither changes the nature of our problem nor imposes substantial penalty in terms of computation complexity.

- p_i can take the form of a lumpsum payment or a stream of per-impression payments. We assume that agents only care about the total expected payment within a given planning horizon.
- p_i can be positive (for buyers) or negative (for the seller). A negative payment means that an agent receives instead of pays money.

We are interested in solving the optimal allocation under each contingency, i.e., for each realized supply vector w , we are interested in x^* that solves

$$(P) \max_x \sum_{i \in \mathcal{N}} u_i(x_i) \quad (5)$$

s.t. (1) and (2)

Combining the optimal allocation for all contingencies, we can obtain the optimal contingent contracts ξ^* .

There are several reasons for us to focus on maximizing total social surplus. First, maximizing social surplus is a good long-term strategy that is useful for maintaining a good long-term business relationship with buyers and for fending off competition. Second, solving the efficient allocation problem is also a key step in designing incentive compatible mechanisms, as we discuss in Section 6.

The following result provides a way of determining whether an allocation x is optimal.

Proposition 1. *An allocation x is optimal if and only if x satisfies constraints (1) and (2), and there exists a shadow price vector $\lambda = (\lambda_1, \lambda_2, \dots, \lambda_M)$ such that*

$$\begin{cases} x_{im} > 0 \Rightarrow \frac{\partial Q_i(x_i)}{\partial x_{im}} = \lambda_m \\ x_{im} = 0 \Rightarrow \frac{\partial Q_i(x_i)}{\partial x_{im}} \leq \lambda_m \end{cases}, \quad \forall i, m \quad (6)$$

λ contains the Lagrange multipliers for constraints (1) in problem (P).

The optimal allocation problem (P) is a nonlinear optimization problem involving $M \times N$ decision variables and $M \times N + M$ constraints. In display advertising both the number of categories M and the number of agents N are large. The former could be in tens of thousands and the latter can be easily in hundreds. Because of the scale of the problem and potentially complex valuation functions, a brutal force attack of the problem (P) may not be feasible.

To address this problem, we first solve the single-category case (i.e., $M = 1$), which serves both as an illustration of our contingent contract approach and as the basis for solving the multi-category case (i.e., $M \geq 2$) in Section 5.

4 Optimal Allocation of A Single Category

With one category, we simplify the notation by replacing x_{im} with x_i , ξ_{im} with ξ_i , α_{im} with α_i , and w_{im} with w_i . We denote

$$V_i \equiv Q'_i(0), \forall i$$

as agent i 's *dropout price*. Intuitively, agent i would drop out if price per impression is higher than V_i . Without loss of generality, we assume

$$V_1 \geq V_2 \geq \dots \geq V_s \geq \dots \geq V_N \geq 0 \quad (7)$$

We denote $d_i(y)$ as agent i 's *induced demand*, i.e., the quantity demanded by agent i if impressions were sold at a price y . Formally,

$$d_i(y) = \begin{cases} Q_i'^{-1}(y), & \text{if } y \leq V_i \\ 0, & \text{otherwise} \end{cases} \quad (8)$$

Clearly, when the price exceeds V_i , agent i demand 0 unit. Also, because the valuation function is strictly concave, the induced demand is monotone decreasing in price.

Proposition 2. *The optimal allocation x^* is uniquely determined via*

$$x_i = d_i(\lambda) \quad (9)$$

where the shadow price $\lambda > 0$ is the unique solution to

$$w = \sum_{i=1}^N d_i(\lambda). \quad (10)$$

To facilitate understanding, we introduce the concept of “floors.” Floors are chunks of total supply ordered by risk levels, with the lowest floor being least risky. Formally, we denote \bar{w}_k as the aggregate demand at agent k 's dropout price; that is,

$$\bar{w}_k = \sum_{i \in \mathcal{N}} d_i(V_k), \quad 1 \leq k \leq N \quad (11)$$

We assume the upper bound of total supply exceeds \bar{w}_N and denote this upper bound as \bar{w}_{N+1} . By construction,

$$0 = \bar{w}_1 \leq \bar{w}_2 \leq \dots \leq \bar{w}_N \leq \bar{w}_{N+1}.$$

We define the k -th floor, f_k , as the chunk of realized supply between \bar{w}_k and \bar{w}_{k+1} , i.e.,

$$f_k = \begin{cases} 0, & w < \bar{w}_k \\ w - \bar{w}_k, & \bar{w}_k \leq w \leq \bar{w}_{k+1} \\ \bar{w}_{k+1} - \bar{w}_k, & w > \bar{w}_{k+1} \end{cases}$$

and we say agent i *participates in* the k -th floor if $\xi_i(w) > 0$ for $w \in (\bar{w}_k, \bar{w}_{k+1}]$. By definition, the total supply can be split into N floors, i.e.,

$$w = \sum_{k=1}^N f_k$$

and higher floors are riskier because they are more likely to be empty.

Proposition 2 implies the following properties about the optimal allocation:

- When the total supply is at the level of k -th floor (i.e., $w \in [\bar{w}_k, \bar{w}_{k+1}]$), the shadow price is between V_{k+1} and V_k .
- The agent k participates only in the k -th floor or above. This is because if the total supply is less than \bar{w}_k , the shadow price is above V_k , agent k 's drop out price. By implication, the seller s retains some impressions if and only if the total supply exceeds \bar{w}_s .
- The quantity offered to each agent is continuous and weakly monotone increasing in the total supply w , as shown by the following corollary.

Corollary 1. *For any $w \in [\bar{w}_k, \bar{w}_{k+1}]$ and $i \leq k$*

$$\frac{d\xi_i^*}{dw} = \frac{\tau_i(x_i)}{\sum_{j=1}^k \tau_j(x_j)} > 0$$

where

$$\tau_i(x) \equiv -\frac{Q'_i(x)}{Q''_i(x)} \quad (12)$$

measures agent i 's absolute risk tolerance.⁴

Corollary 1 suggests that the marginal supply is divided among eligible agents in proportion to their instantaneous absolute risk tolerance: the higher the absolute risk tolerance, the higher the proportion.

⁴The absolute risk tolerance is the reciprocal of the Arrow-Pratt measure of absolute risk aversion.

Algorithm 1 Optimal Single-category Allocation

Input:

- Valuation functions $\{Q_i\}$ and marginal valuation functions $\{Q'_i\}$ in either functional or numerical form.

Steps:

1. Compute dropout prices V_i , $i \in \mathcal{N}$ and reindex agents according to (7).
 2. Compute the quantity thresholds \bar{w}_k , $k \in \mathcal{N}$ which are defined in (11).
 3. For w from 0 to \bar{w} by step Δw do
 - (a) For $w \in [\bar{w}_k, \bar{w}_{k+1}]$, compute λ as the solution to (10) on the interval $[V_{k+1}, V_k]$.
 - (b) Output the optimal allocation x^* as defined by (9)
-

With CARA valuation functions,

$$Q_i(x) = V_i \alpha_i^{-1} (1 - e^{-\alpha_i x}) \quad (13)$$

, all agents have constant risk tolerance

$$\tau_i = \alpha_i^{-1},$$

so agent i owns fixed proportions of the i th floor and above and the optimal allocation is simply,

$$x_i^* = \sum_{k=i}^N \left(\frac{\alpha_i^{-1}}{\sum_{j=1}^k \alpha_j^{-1}} f_k \right) \quad (14)$$

This optimal allocation is extremely easy to implement.

Based on Proposition 2, we develop Algorithm 1 for calculating the optimal contingent contracts under general utility functions. It should be noted that our algorithm does not place any restriction on the form of valuation functions as long as they are strictly increasing and strictly concave. These valuation functions do not have to take any explicit functional form and could be numerically derived from empirical data.

It is straightforward to extend our one-period results to multi-period settings. In particular, if valuation functions and the supply distribution remain the same in each period, we may solve the optimal contingent contracts once and apply the solution repeatedly. If the environment changes, we just have to resolve the optimal contingent contracts.

We conclude this section with numerical simulations of optimal allocation between two agents

$i = 1, 2$ over 100 periods (Figure 2). In these simulations, the total supply is normally distributed with mean 10 and variance 1, and the valuation functions take the CARA form (13).

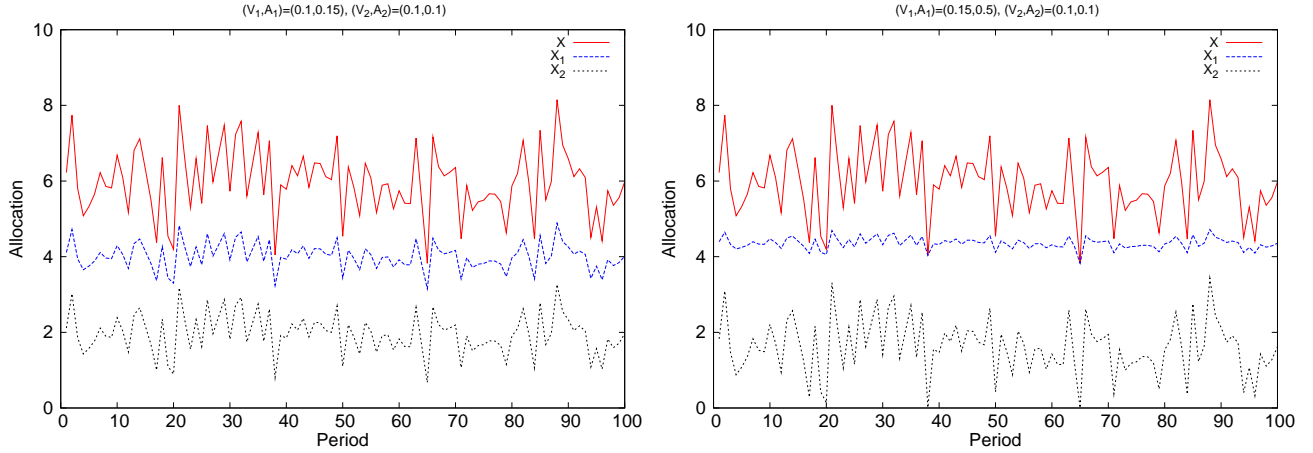


Figure 2: Optimal Allocation Over Time

Figure 2 shows the total impressions (solid red line), and the impressions allocated to agents 1 (dashed blue line) and 2 (dotted black line). In Simulation 2 (a), we let $\mathbf{V} = (0.3, 0.1)$ and $\alpha = (0.5, 0.1)$ so that agent 1 has lower risk tolerance than agent 2 and also higher dropout price. As expected, agent 1's allocation is less volatile. Agent 2's allocation is higher than 1's when the supply is ample but smaller than 1's when the supply is insufficient. In Simulation 2 (b), we let $\mathbf{V} = (0.9, 0.1)$ and $\alpha = (0.5, 0.1)$ so that agent 1's dropout price is even higher. As a result, 2's allocation is sometimes zero and is always less than 1's allocation. In both simulations, agent 1's allocation is less volatile, which is a direct result of her low risk tolerance.

5 Optimal Allocation of Multiple Categories

When the demand for one category of impressions is independent of that for other categories, we can still apply the approach in Section 4 for each category.⁵ Furthermore, if the demand for any two categories are related but the rate of substitution is the same across all agents, we may convert all categories into one and solve the optimal single-category allocation problem. In this section, we focus on the nontrivial case in which demand for different categories are related and the rates of substitutions are not the same.

⁵Note that our approach is applicable even if attributes are revealed in real time, such as in the case of a web user's browsing history. The difference in the realized number of certain impressions amounts to the uncertainty in the supply.

To efficiently solve the multi-category allocation problem, we exploit two special features of the problem. First, a necessary condition for an allocation to be optimal is that it is Pareto optimal. Second, even though agents have nonlinear valuation functions, the marginal rates of substitution between categories are constant. The two special features allow us to decompose the multi-category allocation problem and convert each sub-problem to a standardized single-category allocation problem.

5.1 A Motivating Example

Example 1. We consider the following 4×4 example with CARA valuation functions

$$Q_i(x) = V_i \left(1 - e^{-\sum_{m=1}^4 \alpha_{im} x_{im}} \right), \quad i = 1..4$$

and demand/supply profile

$$\mathbf{V} = \begin{pmatrix} 2 \\ 1 \\ 1.5 \\ 1.2 \end{pmatrix}, \quad \alpha = \begin{bmatrix} 0.30 & 0.16 & 0.10 & 0.20 \\ 0.20 & 0.50 & 0.12 & 0.05 \\ 0.13 & 0.10 & 0.40 & 0.08 \\ 0.06 & 0.10 & 0.20 & 0.30 \end{bmatrix}, \quad w = \begin{pmatrix} 12 \\ 8 \\ 6 \\ 6 \end{pmatrix}$$

We first consider a “naive” approach that allocates each category independently. We obtain the following allocation and social surplus (S).

$$x^0 = \begin{bmatrix} 5.4368 & 5.0289 & 1.3382 & 3.8107 \\ 2.6622 & 2.5018 & 0 & 0 \\ 3.9009 & 0.4693 & 3.0811 & 0 \\ 0 & 0 & 1.5807 & 2.1893 \end{bmatrix}, \quad S = 4.7556.$$

The above allocation is not optimal as there are many opportunities for Pareto improvements. For example, if agent 1 gives one unit of category 2 to agent 2 in exchange of one unit of category 1, then both will be better off and the social surplus will increase to 4.8085. This Pareto-improving trade between agents 1 and 2 can continue until agent 1 runs out of category 2 or agent 2 runs out of category 1. After all such Pareto-improving trades are carried out, we will arrive at a *Pareto optimal allocation*.

But a Pareto optimal allocation is still not necessarily optimal.

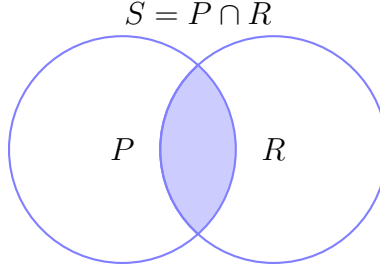


Figure 3: P is the set of allocations that are Pareto optimal and R is the set of allocations that are regular. At least one optimal allocation resides in S .

The optimal allocation and social surplus for this example are

$$x^* = \begin{bmatrix} 11.823 & 0 & 0 & 0 \\ 0 & 6.7119 & 0 & 0 \\ 0.177 & 0 & 6 & 0 \\ 0 & 1.2881 & 0 & 6 \end{bmatrix}, S = 5.3001.$$

As this example shows, joint allocation of multiple categories may result in significant efficiency gains, compared with the naive approach. The optimal allocation matrix above is considerably more sparse because of the Pareto optimality requirement and an additional requirement call regularity. We define these two concepts and study their implications in Sections 5.2 and 5.3 respectively. A crucial result we will obtain, as is illustrated in Figure 1, is that at least one optimal allocation resides in the intersection of the set of Pareto optimal allocations and the set of regular allocations. This result is the foundation of our approach for solving the original problem (P).

5.2 Pareto Optimality

As we have seen in Example 1, Pareto optimality is closely related to pareto-improving trading among agents. Hence, we will define and study Pareto optimality through the notion of trading.

5.2.1 Defining Pareto Optimality

A trade is a reallocation of impressions among agents. A trade T can be represented by a weighted directed graph in which,

- each node is an agent-category pair (i, m) and

- each arc $(i, m) \rightarrow (j, n)$ with weight $\epsilon > 0$ represents a flow of ϵ units of category m from agent i to agent j .⁶

Figure 4 shows a graphical representation of a trade.

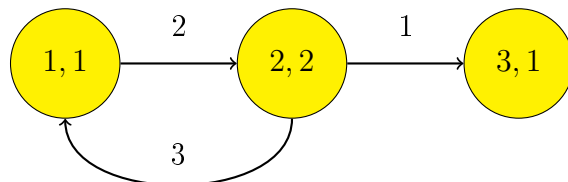


Figure 4: This trade graph represents that agent 1 gives agent 2 two units of category 1; agent 2 gives agent 3 one unit of category 2; agent 2 gives agent 1 three units of category 2.

Let $T(x)$ denote the allocation after a trade T is executed on allocation x .

Definition 1. A trade T is *feasible* on allocation x if and only if $T(x)$ is nonnegative.

Definition 2. A trade T is *profitable* if $u_i(T(x)) \geq u_i(x)$, $\forall i \in \mathcal{N}$ and at least one strict inequality holds.

Definition 3. A trade T is *profit neutral* if $u_i(T(x)) = u_i(x)$, $\forall i \in \mathcal{N}$.

Definition 4. A trade T is *unprofitable* if there exists $i \in \mathcal{N}$ such that $u_i(T(x)) < u_i(x)$.

Definition 5. An allocation x is *Pareto optimal* if none of the feasible trades on x is profitable.

5.2.2 Circular Trades and Trading Cycles

We argue that in order to check Pareto optimality, we only need to look at the structural aspect of a trade and we can focus on a special kind of trades, circular trades. The former is straightforward: if a trade T with weight vector ϵ is feasible (profitable), a scaled version of the trade T with weight vector $\beta\epsilon$ ($\beta \leq 1$) is also feasible (profitable). We will focus on the second argument next.

A *circular trade* is a trade with a circular pattern: i_1 gives ϵ_1 units of category m_1 to i_2 , i_2 gives ϵ_2 units of category m_2 to i_3 , ..., and i_K gives ϵ_K units of category m_K to i_1 . We may concisely denote a circular trade as (C, ϵ) where

$$C = \left((i_1, m_1) \ (i_2, m_2) \ \dots \ (i_K, m_K) \right) \quad (15)$$

⁶We require $i \neq j$ and $m \neq n$ to avoid trivialty. Trade between nodes of the same category can always be avoided because $(i, m) \rightarrow (j, m)$ is equivalent to $(i, m) \rightarrow (j, n)$.

captures the structural aspect of the trade and $\epsilon = \left(\epsilon_1 \ \epsilon_2 \ \dots \ \epsilon_K \right)$ captures the trading quantities. We call C a *trading cycle*.

We say a trading cycle C is feasible (profitable, profit neutral, unprofitable) if there exists a trading quantity vector ϵ such that the circular trade (C, ϵ) is feasible (profitable, profit neutral, unprofitable). Clearly, C is feasible on x if and only if

$$x_{i_k m_k} > 0, \forall k = 1..K \quad (16)$$

This suggests that feasible cycles can only consist of non-zero elements of an allocation matrix. The following results provide a criterion for judging whether a trade cycle is profitable.

Proposition 3. *A trading cycle $C = ((i_1, m_1) \ (i_2, m_2) \ \dots \ (i_K, m_K))$ is profitable (profit neutral, unprofitable) if and only if ⁷*

$$\prod_{k=1}^K \alpha_{i_k m_k} < (=, >) \prod_{k=1}^K \alpha_{i_k m_{k-1}} \quad (17)$$

By reversing the direction of all flows on trading cycle C , we can get a *counter cycle*, denoted as C^{-1} . Intuitively, a circular trade followed by its “counter” circular trade will leave each agent involved in the trade exactly the same.

Example 2. Continue with Example 1. Under the initial allocation x^0 , the trading cycle $C = \left((1, 2) \ (2, 1) \right)$ is both feasible and profitable. The former is because $x_{12} > 0$ and $x_{21} > 0$. The latter is because $\alpha_{12}\alpha_{21} = 0.032 < 0.15 = \alpha_{11}\alpha_{22}$ by Proposition 3.

Proposition 4. *If a trading cycle C is profitable (profit neutral, unprofitable), then its counter cycle C^{-1} is unprofitable (profit neutral, profitable).*

The following proposition implies that we only need to check all trading cycles to know whether an allocation is Pareto optimal.

Proposition 5. *An allocation is Pareto optimal if and only if none of the feasible trading cycles is profitable.*

Example 3. Continue with Example 1. Under the optimal allocation x^* , the trading cycle $C^{-1} = \left((1, 1) \ (1, 2) \right)$ is feasible but not profitable because $\alpha_{11}\alpha_{22} = 0.15 < 0.032 = \alpha_{12}\alpha_{21}$. One can go on to verify that none of the feasible cycles in x^* is profitable.⁸ So, by Proposition 5, x^* is Pareto optimal.

⁷Throughout the paper $i_0, m_0, i_{K+1}, m_{K+1}$ mean i_K, m_K, i_1, m_1 , respectively.

⁸We leave this verification to interested readers. To enumerate all feasible cycles, we may follow three rules: (a) A cycle can only consist of none-zero elements of the allocation matrix (by (16)). (b) Neighbouring nodes should differ in both row numbers and column numbers (see footnote (6)). (3) No two nodes in the cycle should have the same row number (otherwise, we can split the cycle into two smaller cycles).

5.2.3 Pareto Optimal Indicator Matrices and Constrained Problems

Because feasible trading cycles can only consist of non-zero elements of an allocation matrix, it is useful to define the indicator matrix.

Definition 6. A binary matrix D with typical element δ_{im} is an indicator matrix for allocation x if

$$x_{im} = 0 \text{ whenever } \delta_{im} = 0. \quad (18)$$

We also denote $\chi(D)$ as all allocations indicated by D .

It shall be notated that, by our definition, an allocation indicated by D can have either a positive value or zero whenever $\delta_{im} = 1$. For this reason, each D can indicate multiple allocations.

Definition 7. An indicator matrix D is Pareto optimal if all allocations indicated by D are Pareto optimal.

Proposition 6 below provides a sufficient condition for the Pareto-optimality of an indicator matrix.

Proposition 6. An indicator matrix D is Pareto optimal if there exists a strictly positive vector $r = (r_1, r_2, \dots, r_M)$ such that for any $i \in \mathcal{N}$, $m, n \in \mathcal{M}$,

$$\frac{\alpha_{im}}{\alpha_{in}} \geq \frac{r_m}{r_n}, \text{ whenever } \delta_{im} = 1 \quad (19)$$

We interpret r as an *exchange rate vector*. Condition (19) states that agent i can own category m only when her valuation for the category relative to other categories is at least as high as the market's valuation, as defined by the exchange rate vector r . Otherwise, she can profitably trades her category m with others for some other category, as shown in the proof.

Example 4. Continue with Example 1. x^* is indicated by the following binary matrix

$$D = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{bmatrix}.$$

We find that any vector $r = (13, b, 40, 3b)$ with $10 \leq b \leq 20$ satisfies (19). Hence, D is Pareto optimal.

We define a constrained problem with indicator matrix as:

$$\begin{aligned}
 \text{(CP)} \quad \max_x \pi &= \sum_{i \in \mathcal{N}} Q_i(x_i) \\
 \text{s.t.} \quad &(1), (2), \text{ and } x \in \chi(D)
 \end{aligned}$$

We further denote \mathbb{D} as the set of Pareto optimal indicator matrices. Clearly, \mathbb{D} has a finite number of elements.

Because an optimal allocation is necessarily Pareto optimal, the original problem (P) can be solved by solving a series of constrained problems (CP) for each $D \in \mathbb{D}$, until we find one Pareto optimal indicator matrix that yields the optimal allocation for the original problem. The advantage of this approach is the separation of the discrete component (i.e., what categories should be made available to each agent) and the continuous component (i.e., how many units to allocate for each agent and each category) of the problem.

There are three remaining issues:

- First, how do we solve a constrained problem (CP)?
- Second, how do we know whether the solution to a constrained problem (CP) is a solution to the original problem (P)?
- Finally, how do we find an optimal indicator matrix efficiently?

With Proposition 1, the second issue is easy. The remaining two issues can be challenging. To address the first issue, we present a technique for solving the problem (CP) when the optimal indicator matrix is both connected and regular (Section 5.3) and a complementary decomposition method when the optimal indicator matrix is not connected (Section 5.4). To address the last issue, we present some heuristics for matrix search in Section 5.5 (see Appendix B and C for a complete algorithm).

5.3 Regularity

When an agent owns two categories simultaneously, the two categories can be traded at the rate consistent with the agent's marginal rate of substitution, which is determined by her valuation coefficients. When a second agent also owns the same two categories, the two categories can be traded at the second agent's marginal rate of substitution, which causes profitable or profit-neutral trading cycles between the two agents, as shown by the following example.

Example 5. Consider a 2×2 example and an allocation $x = \begin{pmatrix} 2 & 2 \\ 2 & 2 \end{pmatrix}$. If $\alpha_{11}/\alpha_{12} < \alpha_{21}/\alpha_{22}$, then there exists a profitable cycle $C = \left((1,2) \ (2,1) \right)$ by Proposition 3. Thus the allocation is not Pareto optimal. If $\alpha_{11}/\alpha_{12} = \alpha_{21}/\alpha_{22}$, then the cycle C is profit neutral, implying that agents are indifferent between x and allocation $x' = \begin{pmatrix} 2-t & 2+2t \\ 2+t & 2-2t \end{pmatrix}$ for any $0 \leq t \leq 1$. This suggests that if x is optimal so is x' .

If we allow only one agent to own both categories in above example, we can avoid the two kinds of cycles. This requirement is generalized to the multi-category case as the regularity condition. Note though, when there are multiple categories, two categories may be indirectly “connected” via a series of agents instead of directly “connected” via a single agent.

Definition 8. In an indicator matrix D , two categories, m and n , are *connected via agent* i if $\delta_{im} = \delta_{in} = 1$.

Based on the connectivity information in D , we can construct an undirected *connectivity graph* G in which

- each node represents a category, and
- a labeled edge $m \overset{i}{\leftrightarrow} n$ represents m and n are connected through i .

We say D is connected if its corresponding connectivity graph is connected.

Example 6. Figure 5 shows the connectivity graphs corresponding to D^* , D^1 and D^2 where

$$D^* = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{bmatrix}, \quad D^1 = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{bmatrix}, \quad D^2 = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \end{bmatrix}$$

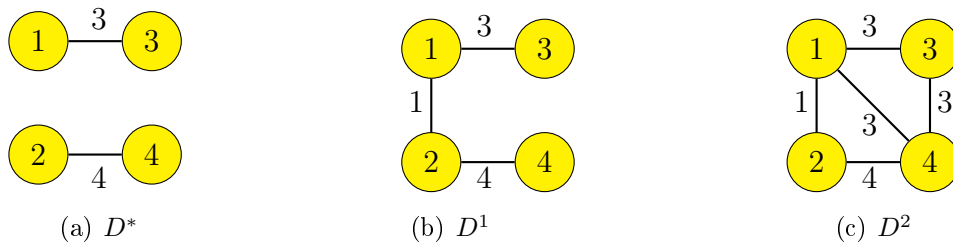


Figure 5: Connectivity graphs corresponding to D^* , D^1 and D^2

Definition 9. A category m is connected to a set of other categories S via agent i if m is connected to at least one category in S via agent i .

Definition 10. The connection between a category m and a set of categories S is *regular* if m is connected to S via a single agent.

In above definition, the connection between a category m and a set S is regular does not necessarily mean that m has a single connection with S . In fact, m can have many connections with S , as long as they are all via the same agent. For example, in D^2 (Figure 5), category 3's connection with $\{1, 2, 4\}$ is regular, as it is via a single agent 3.

Definition 11. A connected component S is *regular* if the connection between any category $m \in S$ to any connected component among the remaining categories is regular.

Note that the remaining categories may form a connected component or several components. The former case occurs with categories 3 and 4 in D^1 (Figure 5) and the latter, with categories 1 and 2. In the latter case, m must be connected to each component via a single agent (though these agents need not to be the same).

Definition 12. An indicator matrix D is *regular* if each of its components is regular.

We denote the set of regular indicator matrices as \mathbb{R} . In Figure 5, D^* is regular because each component of D^* is regular. D^1 is regular. D^2 is not regular because connections between 1, 2, and 4 and the rest are not regular.

Proposition 7. *If a Pareto optimal allocation x is not indicated by any regular indicator matrix, then there must be a profit neutral trade T such that $x' = T(x)$ is indicated by a regular indicator matrix.*

Proposition 7 suggests that we can focus on regular Pareto-optimal indicator matrices without loss of efficiency. We only consider regular indicator matrices from now on. Also, for the purpose of solving (CP), we assume without loss of generality that each row in the indicator matrix D has at least one nonzero element.

Proposition 8. *Consider an allocation problem defined by a connected and regular indicator matrix D that has L connections, $m_1 \overset{i_1}{\leftrightarrow} n_1, m_2 \overset{i_2}{\leftrightarrow} n_2, \dots, m_L \overset{i_L}{\leftrightarrow} n_L$. Then the following L equations determine a vector $r = (r_1, r_2, \dots, r_M)$, which is unique up to a scale factor,*

$$\frac{r_{m_l}}{r_{n_l}} = \frac{\alpha_{i_l m_l}}{\alpha_{i_l n_l}}, \forall l = 1..L. \quad (20)$$

Note that the number of equations L may exceed the number of decision variables. Yet, Proposition 8 suggests that there will still be a unique exchange vector (up to a scale factor). The fact that we can always calculate a unique exchange rate vector for a connected and regular problem is crucial for the following standardization procedure.

Given a connected and regular indicator matrix D and an associated exchange rate vector r as defined by (20), we let

$$\acute{w} = \sum_{m \in \mathcal{M}} w_m r_m \quad (21)$$

be the total number of standardized impressions and

$$\acute{u}_i(\acute{x}_i) = \begin{cases} Q_i \left(\frac{\alpha_{im}}{r_m} \acute{x}_i \right), & \text{if } \exists m, \text{ s.t.}, \delta_{im} = 1, \forall i \in \mathcal{N} \\ 0, & \text{otherwise} \end{cases} \quad (22)$$

be the valuation functions for standardized impressions.⁹

A standardized problem is defined as:

$$(SP) \quad \max_{\acute{x}} \quad \sum_{i \in \mathcal{N}} \acute{u}_i(\acute{x}_i) \\ \text{s.t.} \quad \sum_{i \in \mathcal{N}} \acute{x}_i \leq \acute{w} \text{ and } \acute{x}_i \geq 0, \forall i \in \mathcal{N}$$

The next two results ensure that we can actually recover a solution of the constrained problem (CP) from solving the corresponding standardized problem (SP).

Lemma 1. *If an indicator matrix D of size $N \times M$ is connected and regular, then it has exactly $N + M - 1$ “1” elements.*

Let \acute{x}^* denote the optimal solution to the standardized problem (SP). Let an allocation x be defined by the following system of linear equations:

$$\begin{cases} \sum_{m \in \mathcal{M}} r_m x_{im} \delta_{im} = \acute{x}_i^*, \forall i \in \mathcal{N} \\ \sum_{i \in \mathcal{N}} x_{im} = w_m, \forall m \in \mathcal{M} \end{cases} \quad (23)$$

Proposition 9. *Given a constrained problem (CP) defined by a connected and regular matrix D and the corresponding standardized subproblem (SP). The system of linear equations in (23) always has a unique solution x . Moreover, if x is non-negative, then x is the solution to the constrained problem (CP).*

⁹ $\acute{u}_i(\acute{x})$ is uniquely defined because if $\delta_{im} = \delta_{in} = 1$, we must have $\frac{\alpha_{im}}{r_m} = \frac{\alpha_{in}}{r_n}$ (by Proposition 8).

5.4 Decomposition and Aggregation

We now discuss the case where the indicator matrix that defines a constrained problem (CP) is not connected. Suppose the connectivity graph has J ($1 \leq J \leq M$) connected components. We let \mathcal{M}_j be the nodes in the j th component, \mathcal{N}_j be the set of *affiliated* agents – that is, agents who may hold at least one category in \mathcal{M}_j , and D_j be the submatrix of D defined by rows \mathcal{N}_j and columns \mathcal{M}_j . By construction, each agent is affiliated with at most one component. In this way, we decompose the constrained problem (CP) into J subproblems. In the j -th sub problem, we allocate categories \mathcal{M}_j among agents \mathcal{N}_j , subject to indicator matrix D_j that is connected.

Once we solve each subproblem, we can easily form a candidate solution for the whole problem. To check whether the candidate solution is optimal, we use the following result.

Proposition 10. *Let x be the candidate solution assembled from solving J subproblems and λ_m be the Lagrange multiplier for category $m \in \mathcal{M}_1 \cup \mathcal{M}_2 \dots \cup \mathcal{M}_J$. x is the solution to the original problem (P) if x is non-negative and for each $m \in \mathcal{M}_j$ and $i \notin \mathcal{N}_j$,*

$$\frac{\partial Q_i}{\partial x_{im}} \leq \lambda_m \quad (24)$$

Propositions (9) and (10) together ensure the validity of our approach. We use an example to illustrate the whole process for solving a multi-category allocation problem.

Example 7. Continue with Example (4). We illustrate the decomposition and standardization techniques using the following Pareto optimal indicator matrix,

$$D^* = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{bmatrix}.$$

The first step is to decompose D into two submatrices.

$$D_1 : \begin{array}{cc} & \begin{array}{cc} \text{Category 1} & \text{Category 3} \end{array} \\ \begin{array}{c} \text{Agent 1} \\ \text{Agent 3} \end{array} & \begin{array}{cc} 1 & 0 \\ 1 & 1 \end{array} \end{array}, \quad D_2 : \begin{array}{cc} & \begin{array}{cc} \text{Category 2} & \text{Category 4} \end{array} \\ \begin{array}{c} \text{Agent 2} \\ \text{Agent 4} \end{array} & \begin{array}{cc} 1 & 0 \\ 1 & 1 \end{array} \end{array}$$

The second step is standardization. We use D_2 to illustrate this step. Note that agent 4 has both category 2 and category 4. The exchange rate between the two categories can be calculated as $\alpha_{42}/\alpha_{44} = 1/3$. Let $\begin{pmatrix} r_2 & r_4 \end{pmatrix} = \begin{pmatrix} 1 & 3 \end{pmatrix}$. So the standardized supply

of impressions is $\dot{w} = 8r_2 + 6r_4 = 26$. Agents 2 and 4's utility functions for standardized impressions are

$$\dot{u}_2(\dot{x}_2) = 1 - e^{-0.5\dot{x}_2}, \quad \dot{u}_4(\dot{x}_4) = 1.2(1 - e^{-0.1\dot{x}_4}).$$

The optimal allocation for this single-category problem is

$$\dot{x}_2^* = 6.7119, \quad \dot{x}_4^* = 19.2881.$$

To recover the original allocation, we solve the following system of linear equations.

$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 3 \\ 1 & 1 & 0 \end{pmatrix} \begin{pmatrix} x_{22} \\ x_{42} \\ x_{44} \end{pmatrix} = \begin{pmatrix} 6.7119 \\ 19.2881 \\ 8 \end{pmatrix}$$

where the first two equations are allocative constraints for each agent respectively. The third equation is the allocative constraint for category 2 impressions. The solution to the above system of linear equations is:

$$\begin{pmatrix} x_{22} & x_{24} \\ x_{42} & x_{44} \end{pmatrix} = \begin{pmatrix} 6.7119 & 0 \\ 1.2881 & 6 \end{pmatrix}.$$

For subproblem D_1 , a similar procedure yields

$$\begin{pmatrix} x_{11} & x_{13} \\ x_{31} & x_{33} \end{pmatrix} = \begin{pmatrix} 11.823 & 0 \\ 0.177 & 6 \end{pmatrix}.$$

Both are nonnegative and the combined allocation is

$$x = \begin{bmatrix} 11.823 & 0 & 0 & 0 \\ 0 & 6.7119 & 0 & 0 \\ 0.177 & 0 & 6 & 0 \\ 0 & 1.2881 & 0 & 6 \end{bmatrix}$$

which meets the condition in (24), and thus is optimal. To further check this, we may compute the marginal utility matrix $Q'(x)$ below

$$Q'(x) = \begin{bmatrix} \mathbf{0.017288} & 0.0092202 & 0.0057626 & 0.011525 \\ 0.0069754 & \mathbf{0.017438} & 0.0041852 & 0.0017438 \\ \mathbf{0.017288} & 0.013298 & \mathbf{0.053193} & 0.010639 \\ 0.010463 & \mathbf{0.017438} & 0.034877 & \mathbf{0.052315} \end{bmatrix}.$$

For category- m impression, the value of the Lagrange multiplier λ_m can be found in the m -th column of $Q'(x)$ where the corresponding x is nonnegative (highlighted in bold). The vector of Lagrange multipliers is summarized below:

$$\lambda = (0.017288, 0.017438, 0.053193, 0.052315)'.$$

It's easy to verify that the conditions in Proposition 1 are all satisfied.

5.5 Computational Method

Based on the theory derived in this section, we develop the following algorithm for solving the multi-category optimal allocation problem.

Algorithm 2 Optimal Multi-category Allocation

1. Initialize the indicator matrix $D \in \mathbb{R}$.
 2. Construct the connectivity graph corresponding to D and divide the graph into several connected components.
 3. For each connected component, solve a subproblem by
 - (a) computing the unique exchange rate vector,
 - (b) solving the standardized problem, and
 - (c) backing up the solution for the original subproblem.
 4. Determine whether the candidate solution assembled from solutions of subproblems satisfies the non-negative condition and the optimality condition stated in Proposition 10.
 - (a) If yes, then we have found an optimal allocation.
 - (b) If no, find another $D \in \mathbb{R}$ and go back to Step 2.
-

In Step 1, we may choose any regular Pareto optimal indicator matrix as the initial matrix. For example, one may choose D_0 according to the following rule:

$$\delta_{ij} = \begin{cases} 1 & \text{if } V_i \alpha_{ij} > V_k \alpha_{kj}, \forall k \in \mathcal{N} \\ 0 & \text{otherwise} \end{cases},$$

assuming equality does not occur. Step 2 can be easily achieved using a depth-first search algorithm commonly used in graph theory, and Step 3 can be done using the results in

Section 5.4. There might be many approaches for picking the next D for Step 4 and it is beyond the scope of this paper to evaluate the merits of each approach. Here, we briefly describe the approach we currently use.

Each time after we solve for the constrained problem using the current D matrix, we first check whether there are any negative elements in the allocation matrix. For each negative element x_{im} , we adjust D by setting $\delta_{im} = 0$. After this step, we solve the constrained problem again based on the adjusted D . If all elements in D are nonnegative, then we check the marginal utility of each agent for each category of impression. If for all i and m , the marginal utility $\frac{\partial Q_i}{\partial x_{im}}$ is no larger than λ_m , then we have obtained the optimal allocation. Otherwise, we adjust D by setting $\delta_{im} = 1$ where the ratio $\frac{\partial Q_i}{\partial x_{im}}/\lambda_m$ is the highest. The ratio, which we call the *unbalance* level, captures the increase in marginal valuation, had we reallocate the resource to x_{im} . In choosing the next element to adjust, we also need to consider whether the new D will be still regular. If not, we shall pick the next highest unbalance element, and so on. Once we pick one element to adjust, we solve the constrained problem again based on the adjusted D . This iteration goes on until we find an optimal allocation.

For detailed description of the described algorithm, please see Appendix B. Our preliminary numerical experiments suggest that this heuristic search algorithm can find the optimal allocation very quickly.

Example 8. We illustrate the steps of searching for the optimal D using Example 1. For convenience, we list the demand and the supply parameters below:

$$\mathbf{V} = \begin{pmatrix} 2 \\ 1 \\ 1.5 \\ 1.2 \end{pmatrix}, \alpha = \begin{bmatrix} 0.30 & 0.16 & 0.10 & 0.20 \\ 0.20 & 0.50 & 0.12 & 0.05 \\ 0.13 & 0.10 & 0.40 & 0.08 \\ 0.06 & 0.10 & 0.20 & 0.30 \end{bmatrix}, w = \begin{pmatrix} 12 \\ 8 \\ 6 \\ 6 \end{pmatrix}.$$

1. Set the initial indicator matrix D such that $\delta_{ij} = 1$ if and only if $V_i \alpha_{ij} \geq V_k \alpha_{kj}$, $\forall k \in \mathcal{N}$.

$$D^0 = \begin{bmatrix} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix},$$

By solving the constrained problem defined by D^0 as in Example 7, we have

$$x^0 = \begin{bmatrix} 12 & 0 & 0 & 6 \\ 0 & 8 & 0 & 0 \\ 0 & 0 & 6 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}, \quad Q'(x^0) = \begin{bmatrix} \mathbf{0.0049378} & 0.0026335 & 0.0016459 & \mathbf{0.0032919} \\ 0.0036631 & \mathbf{0.0091578} & 0.0021979 & 0.00091578 \\ 0.017690 & 0.013608 & \mathbf{0.054431} & 0.010886 \\ 0.072000 & 0.12000 & 0.24000 & 0.36000 \end{bmatrix}.$$

2. According to Proposition 10, because some elements of $Q'(x^0)$ exceeds the corresponding λ_m in the same column (shown in bold), x^0 is not optimal. The most unbalanced of all elements is at column 4, row 4, with the ratio $\frac{0.36000}{0.0032919} = 109.36$ (recall the unbalance levels are calculated as elements of $Q'(x^0)$ divided by the Lagrange multiplier, in bold, at the same column). So D^0 is adjusted to

$$D^1 = \begin{bmatrix} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix},$$

and we have

$$x^1 = \begin{bmatrix} 12 & 0 & 0 & -3.3893 \\ 0 & 8 & 0 & 0 \\ 0 & 0 & 6 & 0 \\ 0 & 0 & 0 & 9.3893 \end{bmatrix}, \quad Q'(x^1) = \begin{bmatrix} \mathbf{0.032291} & 0.017222 & 0.010764 & \mathbf{0.021527} \\ 0.0036631 & \mathbf{0.0091578} & 0.0021979 & 0.00091578 \\ 0.017690 & 0.013608 & \mathbf{0.054431} & 0.010886 \\ 0.0043055 & 0.0071758 & 0.014352 & \mathbf{0.021527} \end{bmatrix}$$

3. Because $x_{14} < 0$, we flip δ_{14} and obtain

$$D^2 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix},$$

and we have

$$x^2 = \begin{bmatrix} 12 & 0 & 0 & 0 \\ 0 & 8 & 0 & 0 \\ 0 & 0 & 6 & 0 \\ 0 & 0 & 0 & 6 \end{bmatrix}, \quad Q'(x^2) = \begin{bmatrix} \mathbf{0.016394} & 0.0087436 & 0.0054647 & 0.010929 \\ 0.0036631 & \mathbf{0.0091578} & 0.0021979 & 0.00091578 \\ 0.017690 & 0.013608 & \mathbf{0.054431} & 0.010886 \\ 0.011902 & 0.019836 & 0.039672 & \mathbf{0.059508} \end{bmatrix}$$

4. The most unbalanced element of $Q'(x^2)$ is at column 2, row 4, with ratio $\frac{0.019836}{0.0091578} = 2.166$. So D^2 is adjusted to

$$D^3 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{bmatrix},$$

and we have

$$x^3 = \begin{bmatrix} 12 & 0 & 0 & 0 \\ 0 & 6.7119 & 0 & 0 \\ 0 & 0 & 6 & 0 \\ 0 & 1.2881 & 0 & 6 \end{bmatrix}, \quad Q'(x^3) = \begin{bmatrix} \mathbf{0.016394} & 0.0087436 & 0.0054647 & 0.010929 \\ 0.0069754 & \mathbf{0.017438} & 0.0041852 & 0.0017438 \\ 0.017690 & 0.013608 & \mathbf{0.054431} & 0.010886 \\ 0.010463 & \mathbf{0.017438} & 0.034877 & \mathbf{0.052315} \end{bmatrix}$$

5. The most unbalanced element of $Q'(x^3)$ is at column 1, row 3, with ratio $\frac{0.01769}{0.016394} = 1.08$. So D^3 is adjusted to

$$D^4 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{bmatrix},$$

and we have

$$x^4 = \begin{bmatrix} 11.823 & 0 & 0 & 0 \\ 0 & 6.7119 & 0 & 0 \\ 0.177 & 0 & 6 & 0 \\ 0 & 1.2881 & 0 & 6 \end{bmatrix}, \quad Q'(x^4) = \begin{bmatrix} \mathbf{0.017288} & 0.0092202 & 0.0057626 & 0.011525 \\ 0.0069754 & \mathbf{0.017438} & 0.0041852 & 0.0017438 \\ \mathbf{0.017288} & 0.013298 & \mathbf{0.053193} & 0.010639 \\ 0.010463 & \mathbf{0.017438} & 0.034877 & \mathbf{0.052315} \end{bmatrix},$$

Because x^4 is nonnegative and $Q'(x^4)$ satisfies the optimality condition, the above allocation is optimal and we are done.

For a slightly more complicated example, see Appendix C.

6 Private Information and Incentive Compatible Payments

So far we have examined the optimal allocation problem assuming the planner has complete information about the agents' valuation functions. In this section, we examine the problem

of private information and payments that induce agents to truthfully reveal their preferences.

We assume each agent is characterized by a k -dimensional type vector $t_i = (t_{i1}, t_{i2}, \dots, t_{ik})$. For example, a type vector may include valuation coefficients, reserve valuation, and other parameters of the valuation function. We use t_{-i} and t to denote the types of all but i , and types of all agents respectively.

We denote $u_i(x_i|t_i)$ as the agent i 's valuation of allocation x_i when her type is t_i . We assume u_i is known to the planner but t_i is not.

We apply the well-known VCG mechanism to our settings. In this mechanism, each agent reports his type vector, s_i , to the planner. The planner allocates the impressions by maximizing the total social surplus based on the reported types s and charges each agent $\mu_i(s)$ as defined below.

We denote $x^*(s)$ as the solution to

$$\begin{aligned} \max_x \pi &= \sum_{i \in \mathcal{N}} u_i(x_i|s_i) \\ &s.t. (1) \text{ and } (2) \end{aligned} \quad (25)$$

and $x^{*-i}(s)$ as the solution to

$$\begin{aligned} \max_x \pi &= \sum_{j \in \mathcal{N}} u_j(x_j|s_j) \\ &s.t. (1) \text{ and } (2) \\ &x_i = 0 \end{aligned} \quad (26)$$

In other words, $x^*(s)$ and $x^{*-i}(s)$ are the optimal allocations with and without i respectively.

The VCG mechanism charges agent i a payment of

$$\mu_i(s) = \sum_{j \in \mathcal{N}, j \neq i} u_j(x_j^{*-i}(s)|s_j) - \sum_{j \in \mathcal{N}, j \neq i} u_j(x_j^*(s)|s_j) \quad (27)$$

which is the externality i imposes on other agents by her presence.

Proposition 11. *The multi-dimensional VCG mechanism as defined by the allocation rule (5) and payment rule (27) is truthful in dominant strategies.*

By Krishna and Perry (2000), the multi-dimensional VCG mechanism also maximizes the revenue of the planner among all the efficient mechanisms.

If the seller is the planner, the seller’s type is known and the seller’s total revenue is given by $\sum_{i \neq s} \mu_i(t)$. So the seller’s payoff is

$$U_s(x_s^*(t) | t_s) = u_s(x_s^*(t) | t_s) + \sum_{i \neq s} \mu_i(t)$$

7 Online Allocation and Category Aggregation

We consider the implementation of the optimal allocation in an online environment, where impressions arrive in real time and must be allocated right away. In the “optimize-and-dispatch” framework, our previous discussion is on offline optimization module whereas our next discussion is on online dispatching module. We also discuss the benefits of category aggregation.

7.1 Online allocation

We denote $w(t)$ and $\mathbf{y}(t)$ as the total realized impressions and the total corresponding allocation at discrete time $t \in \{0, 1, \dots, T\}$. The new allocation in period t is thus $\mathbf{y}(t) - \mathbf{y}(t-1)$. We look for an allocation path $\mathbf{y}(t)$ such that

$$\begin{aligned} y_{im}(0) &= 0, \forall i, m, t \\ y_{im}(t) - y_{im}(t-1) &\geq 0, \forall i, m, t \\ \sum_i (y_{im}(t) - y_{im}(t-1)) &\leq w_m(t) - w_m(t-1) \end{aligned}$$

We denote $\xi^*(t) = \{x^{j^*}(t), j = 1, 2, \dots, O\}$ as the set of optimal offline allocations as of time t , i.e., $\xi^*(t)$ is the set of optimal offline allocations if t were the end of the allocation duration. We say $\mathbf{y}(t)$ is an *optimal online allocation* if it satisfies the above conditions and

$$\mathbf{y}(T) \in \xi^*(T).$$

It is easy to see that optimal online allocations are generally not unique, even in the case when the optimal offline allocation is unique.

Next we describe an online allocation algorithm that is close to the offline optimal allocation. We first denote $C_m(t)$ as an *Eligible Set* for a category m at time t . An agent i belongs to $C_m(t)$ if there exists an optimal offline allocation $x^*(t) \in \xi^*(t)$ such that

$$y_{im}(t) < x_{im}^*(t) \tag{28}$$

When an impression of category m comes, we allocate it randomly to any agent in $C_m(t)$. Then we update $C_m(t)$ to $C_m(t+1)$, i.e., remove an agent from $C_m(t)$ if the condition (28) is violated. Intuitively, this process approaches an optimal offline allocation when t reaches T . Note that sometimes multiple impressions of same category may arrive simultaneously. In such a case the consistency of the *Eligible Set* must be preserved to ensure a valid implementation of the optimal allocation. Many existing database techniques can be used for that.

As a further refinement of the online allocation, instead of picking randomly from the Eligible Set, we may also use additional external information to guide the choice. One example of such information is Internet users' preferences, which are often produced by online matching algorithms. By incorporating user preferences, we can better satisfy Internet users while implementing an optimal allocation. More specifically, suppose an online matching algorithm estimates a user's valuation for advertisements as $\{u_i, i = 1..N\}$. Everything else being equal, the seller prefers to assign the impression to the user's most preferred agent. We suppose the seller's concern for advertisers dominates that for Internet users so that the optimal allocation derived in the previous section remains optimal. Thus a real-time refinement can simply be: when an impression of category m arrives, we first find $C_m(t)$ and then pick among $C_m(t)$ that produces the highest user valuation.

7.2 Category aggregation and its implications for computation and privacy

In our main results, all the optimization is done in advance, leaving little to be figured out in real time. In order for such an allocation to be optimal, we need the most detailed demand information in advance, which is not always feasible. Moreover, the global optimization problem can be challenging when the number of categories becomes very large.

To balance between offline and online allocation, we may aggregate the categories to reduce the dimensionality of the global optimization. By doing so, we may lose some information from ad buyers at the offline stage, and the optimal allocation will be in terms of aggregated categories. But at the online stage, we will take into account the detailed preferences within aggregate categories and use them as external information as discussed in Section 7.1.

By reducing the number of categories, the computation cost of the optimization problem can be greatly reduced. At the same time, the size of each eligible set will increase, which means more computation for the online allocation. So this balance of offline optimization and online allocation needs to be tuned carefully to achieve a good tradeoff.

The “optimize-and-dispatch” framework for allocation, together with category aggregation, can also help protect consumer privacy. The online matching, which involves the most detailed Internet user preferences, can be partly carried out in a distributed fashion on Internet user devices such as smart phones. A distributed matching allows Internet users to hold on to their preferences information locally, thus better protecting their privacy. Guha et al. (2011) describes such an online advertising architecture that is considerably more private than current systems, while allowing reasonable ads targeting ads.

With category aggregation, we can limit the advertising providers’ access to the detailed preference information. But the offline optimization based on the coarse information and further refinement by distributed online matching can still provide a reasonable approximation to the optimal allocation. Of course, the computation on user devices should not be too heavy, as the computing power and battery life on smartphones are still very limited. Fortunately there are already some works on these problems in the literature (Guha et al., 2011; Dong et al., 2011). Dong et al. (2011) design a general secure matching algorithm and shows that it is practical, both in terms of computation and energy consumption, with real implementation on smartphones.

8 Summary and Future Research

As more people spend their time online, the demand for efficient allocation of display advertising impressions is unprecedented. We have proposed a contingent contract approach that is flexible enough to accommodate different risk preferences and multiple categories. Although our approach is motivated by the practices of display advertising and is discussed within its context, it is also applicable to sponsored search advertising.

We have shown that with single-category impressions, agents with lower marginal valuation will participate in riskier “floors” of realized impressions. The higher an agent’s absolute risk tolerance, the greater the share he or she will obtain in a floor. An algorithm has been proposed that completely solves the problem of allocating uncertain single-category impressions to multiple agents who have concave valuation functions for impressions.

Our theoretical results on the optimal multi-category allocation shed light on the characteristics of the problem and provides several practical implications. First, our theoretical results ensures the validity of a decomposition and standardization procedure. This procedure suggests a viable algorithm for computing the optimal solution for multiple category allocations. Such a theory-driven algorithm, combined with appropriate matrix search procedure, has potential of achieving superior performance. The byproducts of the decomposition and standardization are also valuable. The result of decomposition informs ad sellers and

buyers which categories of impressions should be considered together. The standardization yields the “exchange rates” across different categories that can further guide buyers in terms of correctly specifying their preferences.

The categories in our framework can capture both the time and location. Therefore, our approach is especially suitable for location- and time-sensitive advertising contexts such as mobile and location-based advertising. The “optimize-and-dispatch” architecture can not only help strike a computational balance between offline optimization and online matching but also lends several additional benefits including better consumer privacy protection. As we have discussed in Section 7, by delegating some matching to online algorithm, one can accommodate consumers’ preference on advertisements, which arrives in real time. Having a distributed online matching can also help preserve consumer privacy, though developing details framework for consumer privacy protection is beyond the scope of this paper.

As a first attack of an inherently challenging problem, we have made several simplifying assumptions. We have assumed that agents know their preferences and assume a VCG type mechanism for soliciting truthful preferences. Other types of preference solicitation mechanisms shall be explored. It will also be valuable to examine in greater detail how to implement optimal allocation in real time. The performance of our approach shall be tested and compared with other approaches in empirical or experimental settings. Despite these assumptions, we believe our existing results show great promise.

References

- Gagan Aggarwal, G. Goel, C. Karande, and A. Mehta. Online vertex-weighted bipartite matching and single-bid budgeted allocations. *CoRR*, 2010.
- Franklin Allen and Douglas Gale. Optimal Security Design. *Review of Financial Studies*, 1(3):229–263, 1989.
- Animesh Animesh, Vandana Ramachandran, and Siva Viswanathan. Quality Uncertainty and the Performance of Online Sponsored Search Markets: An Empirical Investigation. *Information Systems Research*, 21(1):190–201, June 2009.
- S. Athey and G. Ellison. Position auctions with consumer search. *Quarterly Journal of Economics*, 126(3):1213–1270, 2011.
- Jianqing Chen, De Liu, and Andrew B Whinston. Auctioning Keywords in Online Search. *Journal of Marketing*, 73(4):125–141, July 2009.

- Y.J. Chen. Optimal dynamic auctions for display advertising. *UC Berkeley Working Paper*, 2009.
- Wei Dong, Vacha Dave, Lili Qiu, and Yin Zhang. Secure friend discovery in mobile social networks. In *2011 Proceedings IEEE INFOCOM*, pages 1647–1655. IEEE, April 2011.
- Benjamin Edelman, Michael Ostrovsky, and Michael Schwarz. Internet Advertising and the Generalized Second-Price Auction: Selling Billions of Dollars Worth of Keywords. *American Economic Review*, 97(1):242–259, March 2007.
- EMarketer. Online Advertising Market Poised to Grow 20% in 2011. Technical report, 2011.
- A. Ghose and S. Yang. An Empirical Analysis of Search Engine Advertising: Sponsored Search in Electronic Markets. *Management Science*, 55(10):1605–1622, July 2009.
- Arpita Ghosh, P. McAfee, Kishore Papineni, and Sergei Vassilvitskii. Bidding for representative allocations for display advertising. *Internet and Network Economics*, (July 2009): 208–219, 2009.
- Saikat Guha, Bin Cheng, and Paul Francis. Privad: practical privacy in online advertising. In *Proceedings of the 8th Symposium on Networked Systems Design and Implementation (NSDI), Boston, MA*, 2011.
- D Liu and J Chen. Designing online auctions with past performance information. *Decision Support Systems*, 42(3):1307–1320, December 2006.
- D. Liu, J. Chen, and A. B. Whinston. Ex Ante Information and the Design of Keyword Auctions. *Information Systems Research*, 21(1):133–153, November 2010.
- P. McAfee and K. Papineni. Maximally representative allocation for guaranteed delivery advertising campaigns. *Manuscript*, 2010.
- R.P. McAfee. The Design of Advertising Exchanges. *Review of Industrial Organization*, 39: 169–185, 2011.
- Aranyak Mehta, Amin Saberi, Umesh Vazirani, and Vijay Vazirani. AdWords and generalized online matching. *Journal of the ACM*, 54(5):22–es, October 2007.
- D Parkes and T. Sandholm. Optimize-and-dispatch architecture for expressive ad auctions. In *First Workshop on Sponsored Search Auctions, Vancouver, Canada*. Citeseer, 2005.
- A Raviv. The design of an optimal insurance policy. *The American Economic Review*, 69 (1):84–96, 1979.

H Varian. Position auctions. *International Journal of Industrial Organization*, 25(6):1163–1178, December 2007.

R Wilson. Efficient and competitive rationing. *Econometrica*, 57(1):1–40, 1989.

Appendix A

Proof of Proposition 1

Proof. We form the Lagrangian as follows:

$$L = \sum_{i \in \mathcal{N}} Q_i \left(\sum_{m \in \mathcal{M}} \alpha_{im} x_{im} \right) + \sum_{m \in \mathcal{M}} \lambda_m(w) \left(w_m - \sum_{i \in \mathcal{N}} x_{im} \right) + \sum_{m \in \mathcal{M}} \sum_{i \in \mathcal{N}} \mu_i x_{im}.$$

The first order condition for x_{im} is

$$\frac{\partial Q_i}{\partial x_{im}} - \lambda_m(w) + \mu_i = 0$$

Because $\mu_i \geq 0$,

$$\frac{\partial Q_i}{\partial x_{im}} = \lambda_m$$

if $x_{im} > 0$, and

$$\frac{\partial Q_i}{\partial x_{im}} \leq \lambda_m$$

if $x_{im} = 0$. It should be noted that λ_m depends on the vector w . □

Proof of Proposition 2

Proof. We need to show that $\forall k = 1, \dots, N, \forall w \in [\bar{w}_k, \bar{w}_{k+1}]$, the allocation determined by (9) and (10) satisfies the condition in Proposition 1. By construction, the conjectured optimal allocation is feasible provided that the solution to (10) exists. So we first check that (10) indeed has a solution $\lambda(w)$. By the definition of \bar{w}_k and the fact that $w \in [\bar{w}_k, \bar{w}_{k+1}]$, we have

$$\sum_{i=1}^N d_i(V_k) \leq w \leq \sum_{i=1}^N d_i(V_{k+1}).$$

Now because $\sum_{i=1}^N d_i(\cdot)$ is a monotone decreasing function, by continuity, there exists $V_{k+1} \leq \lambda(w) \leq V_k$ such that $w = \sum_{i=1}^N d_i(\lambda(w))$. By Proposition 1, the proposed solution is optimal

if and only if $Q'_i(\xi_i) = \lambda$ holds for $i \leq k$ and $Q'_i(\xi_i) \leq \lambda$ holds for $i > k$. For $i \leq k$, by construction, $Q'_i(\xi_i) = \lambda(w)$, and for $i > k$, $\xi_i = 0$, so

$$Q'_i(\xi_i) = V_i \leq V_k \leq \lambda(w).$$

Therefore, our proposed allocation is indeed optimal. \square

Proof of Corollary 1

Proof. Differentiating $d_i(Q'_i(x_i)) = x_i$ with respect to x_i at $x_i = \xi_i(w)$, we get

$$(Q''_i(x_i))^{-1} = d'_i(Q'_i(x_i)) = d'_i(\lambda(w)).$$

where the second equality is because $Q'_i(x_i) = \lambda(w)$ for $i \leq k$.

Differentiating (9) with respect to w , we get

$$\xi'_i(w) = (Q''_i(x_i))^{-1} \lambda'(w),$$

and, differentiating (10) with respect to w , we get

$$\lambda'(w) = \frac{1}{\sum_{j=1}^k d'_j(\lambda(w))}.$$

Thus,

$$\begin{aligned} \xi'_i(w) &= \frac{(Q''_i(x_i))^{-1}}{\sum_{j=1}^k d'_j(\lambda(w))} = \frac{(Q''_i(x_i))^{-1}}{\sum_{j=1}^k (Q''_j(x_j))^{-1}} \\ &= \frac{(Q''_i(x_i))^{-1} \lambda(w)}{\sum_{j=1}^k (Q''_j(x_j))^{-1} \lambda(w)} = \frac{(Q''_i(x_i))^{-1} Q'_i(x_i)}{\sum_{j=1}^k (Q''_j(x_j))^{-1} Q'_j(x_j)}, \\ &= \frac{R_i(x_i)}{\sum_{j=1}^k R_j(x_j)} \end{aligned}$$

where the fourth step is because $Q'_i(x_i) = \lambda(w)$ for $i \leq k$. \square

Proof of Proposition 3

Proof. (“profitable”) We prove the “if” part by construction. Consider the following circular trade: for each $l = 1, \dots, K$, let agent i_l give $\prod_{k=1}^l \frac{\alpha_{i_k m_{k-1}}}{\alpha_{i_k m_k}}$ units of impression m_l to agent i_{l+1} . The valuation of agent i_l , $l \geq 2$, is given by

$$Q_{i_l} \left(\sum_{m=1}^M \alpha_{i_l m} x_{i_l m} + \alpha_{i_l m_{l-1}} \prod_{k=1}^{l-1} \frac{\alpha_{i_k m_{k-1}}}{\alpha_{i_k m_k}} - \alpha_{i_l m_l} \prod_{k=1}^l \frac{\alpha_{i_k m_{k-1}}}{\alpha_{i_k m_k}} \right) = Q_{i_l} \left(\sum_{m=1}^M \alpha_{i_l m} x_{i_l m} \right)$$

so agent i_l is indifferent after the trade. The valuation of agent i_1 is given by

$$\begin{aligned}
& Q_{i_1} \left(\sum_{m=1}^M \alpha_{i_1 m} x_{i_1 m} + \alpha_{i_1 m_K} \prod_{k=1}^K \frac{\alpha_{i_k m_{k-1}}}{\alpha_{i_k m_k}} - \alpha_{i_1 m_1} \frac{\alpha_{i_1 m_K}}{\alpha_{i_1 m_1}} \right) \\
&= Q_{i_1} \left(\sum_{m=1}^M \alpha_{i_1 m} x_{i_1 m} + \alpha_{i_1 m_K} \left(\prod_{k=1}^K \frac{\alpha_{i_k m_{k-1}}}{\alpha_{i_k m_k}} - 1 \right) \right) \\
&> Q_{i_1} \left(\sum_{m=1}^M \alpha_{i_1 m} x_{i_1 m} \right)
\end{aligned}$$

So agent i_1 is better off after the trade, suggesting a profitable trade on trading cycle C .

We now prove the “only if” part. We suppose a circular trade (C, ϵ) is profitable. Without loss of generality, we assume that i_1 is better off from two adjacent trading steps (i.e., receiving ϵ_K units of m_K from i_K and giving ϵ_1 units of m_1 to i_2) and agents at all other nodes are not worse off from their two adjacent trading steps, namely,

$$\epsilon_K \alpha_{i_1 m_K} > \epsilon_1 \alpha_{i_1 m_1} \quad (29)$$

$$\epsilon_{k-1} \alpha_{i_k m_{k-1}} \geq \epsilon_k \alpha_{i_k m_k}, \forall k = 2..K \quad (30)$$

Multiplying two sides of (29) and (30), we have

$$\prod_{k=1}^K \epsilon_{k-1} \alpha_{i_k m_{k-1}} > \prod_{k=1}^K \epsilon_k \alpha_{i_k m_k}$$

which implies (17).

The proof for the profit-neutral condition is analogous to that for the profitable condition and thus omitted. The condition for unprofitable cycles follows immediately from two previous results. \square

Proof of Proposition 4

Proof. By Proposition 3, the condition for C^{-1} to be profitable (profitable neutral, unprofitable) is

$$\prod_{k=1}^K \alpha_{i_k m_{k-1}} < (=, >) \prod_{k=1}^K \alpha_{i_k m_k} \quad (31)$$

The results in Corollary 4 follows immediately from comparing (17) and (31). \square

Proof of Proposition 5

Proof. We argue that given an allocation x , there is a profitable trade if and only if there is a profitable trading cycle. The Proposition follows naturally from this argument.

The “if” part is obvious. So we only show the “only if” part. First, it is without loss of generality to focus on profitable trades in which each agent both gives and receives. To see, if an agent gives without receiving, the agent is worse off and cannot be part of a profitable trade. If an agent receives without giving, we can drop the agent, return what the agent receives, and obtain a new profitable trade.

Second, given that each agent both gives and receives in the trade t , it always contain a trading cycle. To see, we can start from any agent i in t and trace to someone who receives from i . Because the number of agents is finite, eventually we will reach an agent that we have previously encountered, thus we have a trading cycle C that is feasible under allocation x . If C is profitable, then we have our result. If not, C^{-1} must be profit-neutral or profitable (Proposition 4). So we can find a circular trade (C^{-1}, ϵ) which (a) is profitable or profit neutral and (b) involves each receiver on the trading cycle C returning a portion of the received amount to the sender and at least one receiver returns all the impressions received on C . We can then define a new trade t' that combines t with (C^{-1}, ϵ) . (a) implies that t' is still profitable and (b) implies that t' is still feasible but no longer has the cycle C . (b) also implies that no new trading step, and hence no new cycle, is introduced into t' . Repeating this procedure with t' , there must be another cycle on t' that is either profitable or can be eliminated (without adding new ones) in a new feasible and profitable trade t'' . Doing this repeatedly, eventually, we either find a profitable cycle or there is no cycle left. But the latter is impossible by our earlier argument. \square

Proof of Proposition 6

Proof. We will show that given (19) all feasible cycles are unprofitable. Consider any trading cycle $C = ((i_1, m_1) (i_2, m_2) \dots (i_K, m_K))$. If C is feasible under D , we must have $\delta_{i_k m_k} = 1$ for all $k = 1..K$. By (19), we have

$$\prod_{k=1}^K \frac{\alpha_{i_k m_k}}{\alpha_{i_k m_{k-1}}} \geq \prod_{k=1}^K \frac{r_{m_k}}{r_{m_{k-1}}} = 1$$

which implies that C is not profitable (Proposition 3). \square

Proof of Proposition 7

Proof. We prove by construction. Consider a Pareto optimal allocation x indicated by an irregular D . Suppose that in D , \mathcal{M}_0 is a set of connected categories and $m_0 \notin \mathcal{M}_0$ is connected to $m_1 \in \mathcal{M}_0$ via i_1 and $m_K \in \mathcal{M}_0$ via i_0 ($i_0 \neq i_1$) respectively (the case $K = 1$

is also permitted). Because the categories in \mathcal{M}_0 are connected, there exists a path between m_1 and m_K , say $m_1 \xleftrightarrow{i_2} m_2 \xleftrightarrow{i_3} m_3 \dots \xleftrightarrow{i_K} m_K$ without loss of generality. So the following trading cycle

$$C = ((i_0, m_0) (i_1, m_1) (i_2, m_2) \cdots (i_K, m_K))$$

is feasible. By Proposition 5, C cannot be profitable. If C is unprofitable, then by Proposition 4, C^{-1} is profitable, which cannot be true by Proposition 5 and the fact that C^{-1} is also feasible. Hence, C must be profit neutral. So we can find a profit-neutral trade (C, ϵ) such that after the trade at least one agent i_k runs out of m_k . This is bound to happen because at least two agents are involved in this cycle and $\{m_0, m_1, \dots, m_k\}$ is a distinct set of nodes. If i_0 runs out of m_0 or i_K runs out of m_K , we eliminate a connection between m_0 and \mathcal{M}_0 . If i_k , which is different from i_0 and i_K , runs out of m_k , then either \mathcal{M}_0 is no longer connected, so the premise for the regular connection fails; or we may find a different path connecting m_1 and m_K and repeat the process. The trade does not add new feasible cycles because every recipient of impressions is already allowed to own the category of impressions that he receives. So we can repeatedly use the same technique to eliminate all “redundant” connections or the premise for irregularity without adding new connections or cycles. Because there are only a limited number of redundant connections, we will eventually reach a new allocation x' that is indicated by a regular indicator matrix. \square

Proof of Lemma 1

Proof. We consider the process of constructing D_j step by step. In each step k , an category (column) m_k connected to at least one existing category is added and so are agents (rows) who own m_k but not the existing categories. We denote D_j^k and N_j^k as the indicator matrix and the number of rows after the k th step respectively. Clearly, after adding the first category, D_j^1 has a size of $N_j^1 \times 1$, which has exactly $N_j^1 + 1 - 1$ “1” elements. Suppose D_j^k has $N_j^k + k - 1$ “1” elements. Now we add m_{k+1} . By construction, the new rows contribute exactly $N_j^{k+1} - N_j^k$ “1” elements. By the definition of regular connections, the new column has exactly one “1” element at the existing rows. So the new matrix has $N_j^k + k - 1 + N_j^{k+1} - N_j^k + 1 = N_j^{k+1} + (k + 1) - 1$ “1” elements. By induction, the matrix D_j must have $N_j + M_j - 1$ “1” elements. \square

Proof of Proposition 8

Proof. We show that a vector defined by (20) is unique up to a scale factor. We start by adding a single category m_1 to set \mathcal{M}_0 and let $r_{m_1} = 1$. We then add another category that is connected to \mathcal{M}_0 . By (20), the exchange rate for this new category is uniquely determined. Suppose at the k th step, we have obtained a unique exchange vector $(1, r_{m_2}, r_{m_3}, \dots, r_{m_k})$

according to (20). Now consider adding the $(k + 1)$ th category, m_{k+1} . By regularity, m_{k+1} is connected to \mathcal{M}_0 via a single agent, say i_{k+1} . Suppose m_{k+1} is connected to one or more existing categories, say l_1, l_2, \dots, l_L . If $L = 1$, $r_{m_{k+1}}$ is uniquely determined by (20) as $r_{l_1} \frac{\alpha_{i_{k+1}m_{k+1}}}{\alpha_{i_{k+1}l_1}}$. If $L > 1$, notice that $r_{l_1} \frac{\alpha_{i_{k+1}m_{k+1}}}{\alpha_{i_{k+1}l_1}} = r_{l_2} \frac{\alpha_{i_{k+1}m_{k+1}}}{\alpha_{i_{k+1}l_2}} = \dots = r_{l_L} \frac{\alpha_{i_{k+1}m_{k+1}}}{\alpha_{i_{k+1}l_L}}$, so (20) also uniquely defines $r_{m_{k+1}}$. By induction, the exchange rate vector is well-defined and is unique up to a scale factor. \square

Proof of Proposition (9)

Proof. Because there are exactly $N + M - 1$ “1” elements in the indicator matrix D and there are $N + M - 1$ constraints in (23) after dropping any one redundant constraint implied by (21), so x is exactly determined by (23).

Given x solves (23) and x is nonnegative by assumption, we only need to show that x satisfies 1 for it to be the solution to (CP).

Denote the Lagrange multiplier of the standardized problem as λ . Because \hat{x}^* is the solution to the standardized problem, we have

$$\begin{cases} \hat{x}_i^* > 0 \Rightarrow \left. \frac{\partial \hat{u}_i}{\partial \hat{x}_i} \right|_{\hat{x}_i = \hat{x}_i^*} = \lambda \\ \hat{x}_i^* = 0 \Rightarrow \left. \frac{\partial \hat{u}_i}{\partial \hat{x}_i} \right|_{\hat{x}_i = \hat{x}_i^*} \leq \lambda \end{cases} \quad (32)$$

By Proposition 8, there is an exchange rate vector $r = (r_1 \ r_2 \ \dots \ r_M)$ defined by (20). We show that vector $(\lambda r_1, \lambda r_2, \dots, \lambda r_M)$ satisfies the conditions in (6). First, we notice that if $\delta_{im} = 1$,

$$\sum_{l \in \mathcal{M}} \alpha_{il} x_{il} = \sum_{l \in \mathcal{M}} \alpha_{il} x_{il} \delta_{il} = \frac{\alpha_{im}}{r_m} \sum_{l \in \mathcal{M}} r_l x_{il} \delta_{il} = \frac{\alpha_{im}}{r_m} \hat{x}_i^*$$

where the second equality is because $\frac{\alpha_{im}}{r_m} = \frac{\alpha_{il}}{r_l}$ for any l such that $\delta_{il} = 1$ (see 19). By (22), if $\delta_{im} = 1$,

$$\left. \frac{\partial \hat{u}_i}{\partial \hat{x}} \right|_{\hat{x} = \hat{x}_i^*} = \frac{\alpha_{im}}{r_m} Q'_i \left(\sum_{l \in \mathcal{M}_j} \alpha_{il} x_{il} \right)$$

We discuss the following cases.

(a) If $x_{im} > 0$, then $\hat{x}_i > 0$ and $\delta_{im} = 1$. Hence,

$$\frac{\partial Q_i}{\partial x_{im}} = \alpha_{im} Q'_i \left(\sum_{l \in \mathcal{M}} \alpha_{il} x_{il} \right) = r_m \frac{\alpha_{im}}{r_m} Q'_i \left(\sum_{l \in \mathcal{M}} \alpha_{il} x_{il} \right) = r_m \left. \frac{\partial \hat{u}_i}{\partial \hat{x}_i} \right|_{\hat{x}_i = \hat{x}_i^*} = r_m \lambda$$

(b) If $x_{im} = 0$ and $\delta_{im} = 1$, then, similar to (a), we have

$$\frac{\partial Q_i}{\partial x_{im}} = \alpha_{im} Q'_i \left(\sum_{l \in \mathcal{M}} \alpha_{il} x_{il} \right) = r_m \frac{\alpha_{im}}{r_m} Q'_i \left(\sum_{l \in \mathcal{M}} \alpha_{il} x_{il} \right) = r_m \frac{\partial \acute{u}_i}{\partial \acute{x}_i} \Big|_{\acute{x}=\acute{x}_i^*} \leq r_m \acute{\lambda}$$

where the inequality is due to (32).

(c) If $x_{im} = 0$ and $\delta_{im} = 0$, then there must exist $n \in \mathcal{M}$ and $n \neq m$ such that $\delta_{in} = 1$.

We have

$$\frac{\partial Q_i}{\partial x_{im}} = \frac{\alpha_{im}}{\alpha_{in}} \frac{\partial Q_i}{\partial x_{in}} \leq \frac{\alpha_{im}}{\alpha_{in}} r_n \acute{\lambda} \leq r_m \acute{\lambda}$$

where the first inequality can be seen from steps (a) and (b) and the second inequality is inferred from $\delta_{in} = 1$ and Propositions 6 and 8.

Taking (a), (b), and (c) together, we know x satisfies conditions (6). Because x satisfies the feasibility conditions by construction, if x satisfies the non-negative conditions, x is the solution to the original problem (by Proposition 1) \square

Proof of Proposition 10

Proof. By construction, x satisfies the feasibility condition. With the condition that x is non-negative, we only need to show that x satisfies (6). We have already shown that the conditions (6) are satisfied for (i, m) such that $m \in \mathcal{M}_j$ and $i \in \mathcal{N}_j$. Because for $m \in \mathcal{M}_j$ and $i \notin \mathcal{N}_j$, $x_{im} = 0$, we need condition (24) (note that the Lagrange multipliers are determined by the sub problems). \square

Proof of Proposition (11)

Proof. Notice that the first term is not a function of the agent i 's report s_i , by well-known results (e.g., Proposition 23.C.4 of text book), reporting truthfully is a weakly dominant strategy.

With truthful reporting,

$$\begin{aligned} U_i(x_i^*(t) | t_i) &= u_i(x_i^*(t) | t_i) - \left[\sum_{j \in \mathcal{N}, j \neq i} u_j(x_j^{*-i}(t) | t_j) - \sum_{j \in \mathcal{N}, j \neq i} u_j(x_j^*(t) | t_j) \right] \\ &= \sum_{j \in \mathcal{N}} u_j(x_j^*(t) | t_j) - \sum_{j \in \mathcal{N}, j \neq i} u_j(x_j^{*-i}(t) | t_j) \\ &= \sum_{j \in \mathcal{N}} u_j(x_j^*(t) | t_j) - \sum_{j \in \mathcal{N}} u_j(x_j^{*-i}(t) | t_j) + u_i(0 | t_i) \\ &\geq u_i(0 | t_i) \end{aligned}$$

so the individual rationality constraint for each agent is also satisfied. \square

Appendix B

Before we present the flowchart of the searching algorithm, we need a few definitions. Suppose at step t , we have the allocation matrix x , the associated indicator matrix D . Denote the marginal utility matrixes as follows:

$$Q'(x) = \begin{bmatrix} q_{11} & q_{12} & \cdots & q_{1m} \\ q_{21} & q_{22} & \cdots & q_{2m} \\ \vdots & \vdots & \vdots & \vdots \\ q_{n1} & q_{n2} & \cdots & q_{nm} \end{bmatrix}$$

Define the Lagrange multiplier vector at step t as $\lambda = (\lambda_1, \lambda_2, \dots, \lambda_m)$ where

$$\lambda_j = \max_i q_{ij} \delta_{ij}.$$

Define the balance matrix *unbalance* with typical element b_{ij} defined as

$$b_{ij} = q_{ij} / \lambda_j.$$

Let *unbalancedID* be a list with element *unbalancedID* _{k} being the index of the k -th largest element in *unbalance* matrix that exceeds 1 (elements in *unbalance* are indexed as stacked columns). So *unbalancedID*[0] is most unbalanced element, which will be considered first for adjustment. Also let *unbalancedN* be the size of *unbalancedID*, i.e, the number of greater-than-one elements in *unbalanced*.

Figure 6 shows the flowchart of our preliminary searching algorithm.

Appendix C

The following is a slightly more complicated 5×10 example with CARA valuation functions

$$Q_i(x) = V_i \left(1 - e^{-\sum_{m=1}^{10} \alpha_{im} x_{im}} \right)$$

where the vector \mathbf{V} , the matrix α , and the supply are

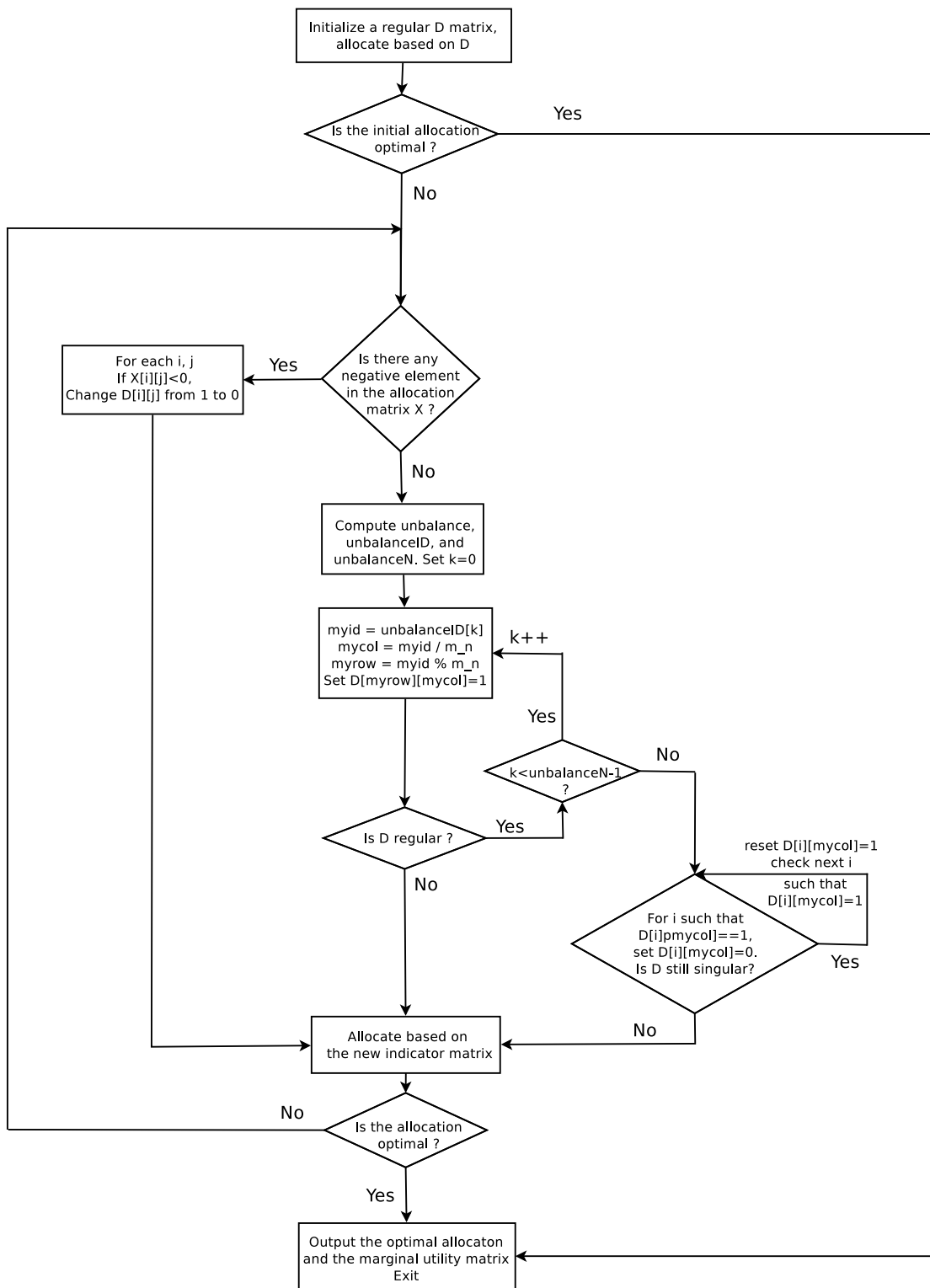


Figure 6: A Searching Algorithm

$$\mathbf{V} = \begin{pmatrix} 2 \\ 1 \\ 1.5 \\ 1.2 \\ 1 \end{pmatrix}, \alpha = \begin{bmatrix} 0.3 & 0.16 & 0.1 & 0.2 & 0.05 & 0.08 & 0.12 & 0.22 & 0.09 & 0.17 \\ 0.2 & 0.5 & 0.12 & 0.05 & 0.1 & 0.11 & 0.2 & 0.08 & 0.07 & 0.02 \\ 0.13 & 0.1 & 0.4 & 0.08 & 0.16 & 0.15 & 0.18 & 0.05 & 0.2 & 0.11 \\ 0.06 & 0.1 & 0.2 & 0.3 & 0.02 & 0.12 & 0.08 & 0.16 & 0.23 & 0.27 \\ 0.16 & 0.11 & 0.05 & 0.2 & 0.06 & 0.03 & 0.13 & 0.17 & 0.09 & 0.21 \end{bmatrix},$$

$$w = (15, 15, 15, 15, 15, 15, 15, 15, 15, 15)'$$

The optimal D matrix and the corresponding optimal allocation are

$$D = \begin{bmatrix} 1 & 0 & 0 & 0 & 1 & 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 1 & 1 & 0 & 1 \end{bmatrix}, x = \begin{bmatrix} 15 & 0 & 0 & 0 & 15 & 0 & 0 & 5.5792 & 0 & 0 \\ 0 & 15 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 15 & 0 & 0 & 13.737 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 15 & 0 & 1.263 & 0 & 0 & 15 & 0 \\ 0 & 0 & 0 & 15 & 0 & 0 & 15 & 9.4208 & 0 & 15 \end{bmatrix}$$

The marginal utility matrix is

$$Q'(x) = \begin{bmatrix} \mathbf{0.00020587} & 0.00010980 & 6.8625e-05 & 0.00013725 & \mathbf{0.00010294} \\ 0.00011062 & \mathbf{0.00027654} & 6.6370e-05 & 2.7654e-05 & 5.5308e-05 \\ 6.1574e-05 & 4.7365e-05 & \mathbf{0.00018946} & 3.7892e-05 & 2.8419e-05 \\ 8.8808e-05 & 5.9206e-05 & 0.00011841 & \mathbf{0.00017762} & 1.1841e-05 \\ 0.00014209 & 9.7689e-05 & 4.4404e-05 & \mathbf{0.00017762} & 5.3285e-05 \end{bmatrix}$$

$$\sim \begin{bmatrix} 5.4900e-05 & 8.2350e-05 & \mathbf{0.00015097} & 6.1762e-05 & 0.00011666 \\ 6.0839e-05 & 0.00011062 & 4.4247e-05 & 3.8716e-05 & 1.1062e-05 \\ \mathbf{7.1047e-05} & 8.5256e-05 & 2.3682e-05 & 9.4729e-05 & 5.2101e-05 \\ \mathbf{7.1047e-05} & 4.7365e-05 & 9.4729e-05 & \mathbf{0.00013617} & 0.00015986 \\ 2.6643e-05 & \mathbf{0.00011545} & \mathbf{0.00015097} & 7.9928e-05 & \mathbf{0.00018650} \end{bmatrix}$$