

Budgeted Online Assignment in Crowdsourcing Markets: Theory and Practice

(Extended Abstract)

Pan Xu, Aravind Srinivasan
Dept. of Computer Science
University of Maryland, College Park, USA
{panxu, srin}@cs.umd.edu

Kanthi K. Sarpatwar, Kun-Lung Wu
IBM Thomas J. Watson Research Center
Yorktown Heights, NY, USA
{sarpatwa, klwu}@us.ibm.com

ABSTRACT

We consider the following budgeted online assignment (BOA) problem motivated by crowdsourcing. We are given a set of offline tasks that need to be assigned to workers who come online from the pool of types $\{1, 2, \dots, n\}$. For a given time horizon $\{1, 2, \dots, T\}$, at each instant of time t , a worker j arrives from the pool in accordance with a known probability distribution $[p_{jt}]$ such that $\sum_j p_{jt} \leq 1$; j has a known subset $N(j)$ of the tasks that it can complete, and an assignment of one task i to j (if we choose to do so) should be done before task i 's deadline. The assignment $e = (i, j)$ (of task $i \in N(j)$ to worker j) yields a profit w_e to the crowdsourcing provider and requires different quantities of K distinct resources, as specified by a cost vector $\mathbf{a}_e \in [0, 1]^K$; these resources could be client-centric (such as their budget) or worker-centric (e.g., a driver's limitation on the total distance traveled or number of hours worked in a period). The goal is to design an online-assignment policy such that the total expected profit is maximized subject to the budget and deadline constraints.

We propose and analyze two simple linear programming (LP)-based algorithms and achieve a competitive ratio of nearly $1/(\ell + 1)$, where ℓ is an upper bound on the number of non-zero elements in any \mathbf{a}_e . This is nearly optimal among all LP-based approaches.

1. INTRODUCTION

Crowdsourcing markets (e.g., Amazon Mechanical Turk or Crowdflower) have evolved to be powerful platforms that bring together task performers (or workers) and task requesters (or consumers). In recent years, problems arising from online decision making in such settings have been attracting tremendous attention (see the survey [5]). A typical problem arising in such settings, considered by [3], is to schedule a batch of consumer tasks using a pool of workers who become available in an online fashion (i.e., in real time). More specifically, we are given a set I of *offline* tasks, where each task $i \in I$ has a deadline d_i after which it cannot be scheduled. Workers arrive in an online fashion (according to an adversarial or random permutation order) and submit bids on a subset of tasks that interest them. When a worker

j arrives, a decision must be made immediately and irrevocably - either assign it an available task or reject its service. If the worker j is allocated a task i , we must pay the worker their bid amount b_{ij} . The goal is to maximize the number of tasks assigned while constrained by a given bid budget of B .

Let $[n]$ denote the set of integers $\{1, 2, \dots, n\}$ and assume a time horizon $[T]$. Our work deals with a natural variant of this problem in the following ways. First we model the arrival of workers as each time $t \in [T]$ a single worker is chosen from a known pool of worker types $[n]$ in accordance with a *known* probability distribution $[p_{jt}]$ such that $\sum_j p_{jt} \leq 1$. Here we allow that with probability $1 - \sum_j p_{jt}$, none of the workers is chosen at t . Second, we consider multiple budget constraints. That is, we assume that there are K distinct resources and that each assignment $e = (i, j)$ has a bid cost vector $\mathbf{a}_e \in [0, 1]^K$, where the k^{th} component of the vector corresponds to the amount of resource type k needed by the assignment. Finally, instead of maximizing the throughput (i.e., number of tasks completed), each assignment e is associated with a *known* weight or utility w_e and we aim to maximize the expected utility collected from those successful assignments.

2. PROBLEM STATEMENT

In this section, we present a formal statement of our problem. Let $I = \{i \in [m]\}$ be the set of offline tasks and $J = \{j \in [n]\}$ be the set of online workers. On a finite time horizon T , each task i has a deadline $d_i \in [T]$ after which it will become unavailable. Let $G = (I, J, E)$ be the bipartite graph that models the relation between the tasks and workers: there is an edge $e = (i, j)$ iff worker j is interested in the task i . Let $N(j) = \{i : (i, j) \in E\}$ be the set of tasks that interest worker j and $N(i) = \{j : (i, j) \in E\}$ be the set of workers who are interested in task i . Each edge $e = (i, j)$ has a weight w_e denoting the profit obtained by assigning task i to worker j . Each assignment $e = (i, j)$ has a requirement for one or more of a given set of K types of resources. The requirement of an assignment e is given by a K -dimensional vector $\mathbf{a}_e \in [0, 1]^K$, where the k^{th} dimension $a_{e,k}$ represents the amount of resource k needed. Each resource type k has a budget $B_k \in \mathbb{R}_+$ that must not be violated. For each e , let $S_e = \{k \in [K] : a_{e,k} > 0\}$, i.e., the set of resources it requests.

Let $E_{j,t} = \{e = (i, j), i \in N(j) : d_i \geq t\}$ denote the set of *available* assignments for the worker j at time t . In this paper, we assume without loss of generality that each

Appears in: *Proc. of the 16th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2017)*, S. Das, E. Durfee, K. Larson, M. Winikoff (eds.), May 8–12, 2017, São Paulo, Brazil.
Copyright © 2017, International Foundation for Autonomous Agents and Multiagent Systems (www.ifaamas.org). All rights reserved.

task can be assigned for an arbitrary number of times before its deadline. Any potential restriction on the number of assignments can easily be modeled by an additional budget constraint: the task itself is an integral resource and the corresponding budget is the upper bound on the number of assignments. For each $e \in E_{j,t}$, we say e is *safe* or *valid* iff for each $k \in S_e$, resource k has remaining budget larger or equal to $a_{e,k}$. When a worker j arrives at t , we have to make an immediate and irrevocable decision: either reject it or choose a safe option $e \in E_{j,t}$ and get a resultant profit w_e . Once a safe assignment e is scheduled, the budget of each resource $k \in S_e$ will be reduced by $a_{e,k}$. Our goal is to design an online assignment policy such that the expected profit is maximized.

In most applications, we need to deal with two kinds of resources, namely integral and non-integral. A resource k is integral if $a_{e,k} \in \{0, 1\}$ for all $e \in E$ and $B_k \in \mathbb{Z}_+$. On the other hand a resource k is non-integral if $a_{e,k} \in [0, 1]$ and $B_k \in \mathbb{R}_+$. This captures resources such as money and time that cannot be quantified as integral. Let $\mathcal{K}_1 = \{1, 2, \dots, K_1\}$ and $\mathcal{K}_2 = \{K_1 + 1, \dots, K_1 + K_2\}$ denote the set of integral and non-integral resources respectively. As defined in the introduction, for each assignment e , $|S_e \cap \mathcal{K}_1| \leq \ell_1$ and $|S_e \cap \mathcal{K}_2| \leq \ell_2$.

3. BENCHMARK LP

Recall that $E_{j,t}$ is the set of available assignments for a worker j arriving at t . For any t , let $E_t = \bigcup_j E_{j,t}$ be the set of all available assignments at t . Further, for each t and $e \in E_t$, let $x_{e,t}$ be the probability that we make the assignment e at t in the offline optimal solution. Our benchmark LP can now be described as follows:

$$\text{maximize } \sum_t \sum_{e \in E_t} w_e x_{e,t} \quad (1)$$

$$\text{subject to } \sum_{e \in E_{j,t}} x_{e,t} \leq p_{jt} \quad \forall j \in J, t \in [T] \quad (2)$$

$$\sum_t \sum_{e \in E_t} x_{e,t} a_{e,k} \leq B_k \quad \forall k \in [K] \quad (3)$$

$$0 \leq x_{e,t} \leq 1 \quad \forall e \in E, t \in [T] \quad (4)$$

LEMMA 1. *The optimal value to LP (1) is a valid upper bound for the offline optimal.*

Our benchmark LP is essentially the same as that used in [1] and [2]. The detailed proof can be found there.

4. TWO LP-BASED ALGORITHMS

In the section, we present two simple LP-based algorithms, ALG_1 and ALG_2 , which are non-adaptive and adaptive respectively. Let $\{x_{e,t}^* | t \in [T], e \in E_t\}$ be an optimal solution for the LP (1).

ALGORITHM 1: A simple non-adaptive algorithm (ALG_1)

For each time t , assume some worker j arrives.

Let $\hat{E}_{j,t} \subseteq E_{j,t}$ be the set of *safe* available assignments we can make for j .

If $\hat{E}_{j,t} = \emptyset$, then reject j ; otherwise sample at most one assignment $e \in \hat{E}_{j,t}$ with probability $\alpha x_{e,t}^*/p_{jt}$.

ALGORITHM 2: Simulation-based adaptive algorithm (ALG_2)

For each time t , assume some worker j arrives.

Let $\hat{E}_{j,t} \subseteq E_{j,t}$ be the set of *safe* available assignments we can make for j .

If $\hat{E}_{j,t} = \emptyset$, then reject j ; otherwise sample an assignment

$e \in \hat{E}_{j,t}$ with probability $\frac{x_{e,t}^*}{p_{j,t}} \frac{\gamma}{\beta_{e,t}}$, where $\beta_{e,t}$ is an estimation of the probability that e is safe at t through simulation.

First, we consider the simple case where all the resources are integral and each assignment requests at most $\ell_1 = \ell$ resources.

THEOREM 1. *By choosing $\alpha = 1/(\ell + 1)$, ALG_1 achieves a competitive ratio of $\frac{1}{\ell+1} (1 - \frac{1}{\ell+1})^\ell \geq \frac{1}{e(\ell+1)}$ for the BOA problem when all the resources are integral.*

THEOREM 2. *By choosing $\gamma = 1/(\ell + 1)$, ALG_2 can achieve a competitive ratio of $(1 - \epsilon)/(\ell + 1)$ for the BOA problem when all the resources are integral for any given $\epsilon > 0$.*

Note that our competitive ratio analysis is tight for the non-adaptive algorithm ALG_1 . Furthermore, ALG_2 is nearly optimal among all approaches based on LP (1) since it has an integrality gap at least $\ell - 1 + 1/\ell$ ([4]).

Second, we consider the general case when both of integral and non-integral resources are involved. Let B be the minimum budget for any non-integral resource.

THEOREM 3. *For the BOA problem, ALG_1 yields a competitive ratio of $\frac{1}{\ell_1+1} \left((1 - \frac{1}{\ell_1+1})^{\ell_1} - \epsilon \right)$, for any $\epsilon > 0$, assuming $B \geq 2 \ln(\frac{\ell_2}{\epsilon}) \left(1 + \frac{3\ell_1+2}{\ell_1^2} \right) + 2$.*

THEOREM 4. *For the BOA problem, ALG_2 yields a competitive ratio of $\frac{1-2\epsilon}{\ell_1+1}$ for any given $\epsilon > 0$, assuming $B \geq 3 \ln(\frac{\ell_2}{\epsilon}) (1 + \frac{1}{\ell_1}) + 2$.*

Our results show that the knowledge about arrival distributions holds a significant edge over the adversarial model or the random permutation model. Let us compare our results with those of [3]. As discussed before, their setting fits our model when $\ell_1 = \ell_2 = 1$. From Theorem 4, we obtain a $(\frac{1}{2} - \epsilon)$ competitive ratio assuming $B \geq 12 \ln(1/\epsilon)$ while [3] obtain a ratio of $O(\frac{1}{R^\epsilon \ln R})$, assuming $B \geq \frac{R}{\epsilon}$ and $R \doteq \frac{\max b_{i,j}}{\min b_{i,j}}$ (i.e., the ratio of the largest bid to the smallest bid over all possible assignments). In fact, we completely remove the dependence on R and obtain a constant ratio while relaxing the lower bound assumption on B significantly. Our result may be seen as theoretical evidence to advocate the use of historical data to learn arrival distributions.

Acknowledgments

Part of this work is done when Pan Xu interned at the IBM T. J. Watson Research Center during the summer of 2016. Aravind Srinivasan's research was supported in part by NSF Awards CNS-1010789 and CCF-1422569, and by a research award from Adobe, Inc. Pan Xu's research was supported in part by NSF Awards CNS-1010789 and CCF-1422569.

REFERENCES

- [1] S. Alaei, M. Hajiaghayi, and V. Liaghat. Online prophet-inequality matching with applications to ad allocation. In *Proceedings of the 13th ACM Conference on Electronic Commerce*, pages 18–35. ACM, 2012.
- [2] S. Alaei, M. Hajiaghayi, and V. Liaghat. The online stochastic generalized assignment problem. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*, pages 11–25. Springer, 2013.
- [3] S. Assadi, J. Hsu, and S. Jabbari. Online assignment of heterogeneous tasks in crowdsourcing markets. In *Third AAAI Conference on Human Computation and Crowdsourcing*, 2015.
- [4] Z. Füredi, J. Kahn, and P. D. Seymour. On the fractional matching polytope of a hypergraph. *Combinatorica*, 13(2):167–180, 1993.
- [5] A. Slivkins and J. W. Vaughan. Online decision making in crowdsourcing markets: Theoretical challenges. *ACM SIGecom Exchanges*, 12(2):4–23, 2014.