

Stochastic Budget Optimization in Internet Advertising

Bhaskar DasGupta
Department of Computer Science
University of Illinois at Chicago
Chicago, IL 60607
dasgupta@cs.uic.edu

S. Muthukrishnan
Department of Computer Science
Rutgers University
Piscataway, NJ 08854
muthu@cs.rutgers.edu

April 4, 2011

Abstract

Internet advertising is a sophisticated game in which the many advertisers “play” to optimize their return on investment. There are many “targets” for the advertisements, and in each “target” is a collection of games with a potentially different set of players are involved. In this paper, we study the problem of how advertisers allocate their budget across these “targets”. In particular, we focus on formulating their *best response* strategy as an optimization problem. Advertisers have a set of keywords (“targets”) and some stochastic information about the future, namely a probability distribution over *scenarios* of cost vs click combinations. This summarizes the potential states of the world assuming that the strategies of other players are fixed. Then, the best response can be abstracted as *stochastic budget optimization* problems to figure out how to spread a given budget across these keywords to maximize the expected number of clicks.

We present the *first known* non-trivial poly-logarithmic approximation for these problems as well as the first known hardness results of getting better than logarithmic approximation ratios in the various parameters involved. We also identify several special cases of these problems of practical interest, such as with fixed number of scenarios or with polynomial-sized parameters related to cost, which are solvable either in polynomial time or with improved approximation ratios. Stochastic budget optimization with scenarios has sophisticated technical structure. Our approximation and hardness results come from relating these problems to a special type of (0/1, bipartite) quadratic programs inherent in them. Our research answers some open problems raised by the authors in (Stochastic Models for Budget Optimization in Search-Based Advertising, Algorithmica, 58 (4), 1022-1044, 2010).

1 Introduction

This paper deals with the problem of how advertisers allocate their budget in Internet advertising. In sponsored search, users who pose queries to internet search engines are not only provided search results, but also a small set of text ads. These ads are chosen from a set of campaigns set up by advertisers based on the keywords in the search query. A lot of focus has been on how these ads are chosen and priced, which is via an auction that is by now well known [2, 9, 23]¹. Our focus is instead on the problem faced by advertisers. Even small advertisers have many *keywords*, a budget in mind and must figure out how to spread this budget on bids for each of these keywords. This is a highly nontrivial task, and the basis for a separate industry to support advertisers. A similar problem arises with “display ads” where advertisers have *websites* where their ads will be shown

¹Likewise, there was a lot of work on bidding strategies [3, 12, 21, 26]. This paper extends that body of work by considering a richer model of uncertainty; see subsequent paragraphs.

and need to split their budget for the ad campaign across the sites to be most effective. Likewise, in behavioral targeting, advertisers have to decide how to spread their budget across *behavior groups*. In all these cases, therefore, advertisers have various “targets” and wish to split their budget across them to optimize their ad campaigns.

Consider the sponsored search example and fix an advertiser A . They have many keywords that they would like to target for their ads. How should they bid for each, given some overall budget they can spend? There is a sophisticated underlying game in which the many advertisers “play” to optimize their return on investment simultaneously. For each keyword and for each instance of auction triggered by this keyword, there is potentially a different set of competing advertiser involved. Building effective strategies is challenging amidst so many parameters. A fundamental and widely accepted proposal is for the advertiser A to pursue a *best response* strategy, *i.e.*, fix the strategies of other advertisers and pick the best strategy as one’s response. Besides being a simple and easy strategy to understand and hence suitable for experimentation by advertisers, best response has desirable properties. For example, in the absence of budgets and for single repeated auction, special type of best response by every player leads to the VCG outcome [4, 5, 9, 23].

In order to help the advertisers implement this best response strategy, search engines provide them with expected bid versus clicks function for each keyword². Assuming that the rest of the world is fixed, these functions provide an estimate of the expected number of clicks an advertiser would obtain by bidding a certain value on that keyword. These functions can also be “learned” by an advertiser to some extent by systematically trying out various bids. Finding advertiser’s best response bidding strategy then becomes an optimization problem where the goal is to maximize the expected number of clicks assuming access to these functions. The resulting problems are in the spirit of the Knapsack problem [3, 12, 21, 26] with many of them solvable nearly exactly or with constant factor approximations³.

A more general approach is to acknowledge that, in reality, the bids vs clicks functions are *not fixed*, but rather random variables with unknown correlations and uncertainties: number of queries (and hence, clicks and budget spent on a keyword) change each day, relative occurrences of keywords change (*e.g.*, searches for beach and snow are complementary⁴), and so on. Therefore, one has to consider a specific stochastic model for these random variables and then maximize the *expected number of clicks* under that model. This approach was initiated in [21] leading to a stochastic budget optimization problem that is studied in this paper.

1.1 Organization of the paper

For convenience of the readers, we organize the rest of the paper in the following manner.

- We start with Section 2 which describes all of our stochastic budget optimization models and corresponding computational problems precisely, starting from the simplest one, together with some comments and justifications about the model. In the last subsection of this section (Section 2.5), we fix some notational uniformity for readers convenience.
- In Section 3, we summarize the results obtained in this paper. For the benefit of the reader, we group the results into two categories, namely a set of main results that deal with the com-

²See, for example, Traffic Estimator at <http://adwords.google.com/support/aw/bin/answer.py?hl=en&answer=8692>, bidding tutorial at <http://adwords.google.com/support/aw/bin/answer.py?hl=en&answer=163828> and bid simulator at <http://adwords.google.com/support/aw/bin/answer.py?hl=en&answer=138148>

³The scenario model was introduced in [21]. For a very detailed discussion of prior works related to the approach in the model, see Section 1.4 of [21].

⁴See www.google.com/trends?q=beach%2C+snow&ctab=0&geo=all&date=all&sort=0 for yearly and www.google.com/trends?q=clubs%2C+stocks&ctab=0&geo=all&date=mtd&sort=0 for weekly trends.

putational complexity issues of the original models without restrictions and a set of additional results that deal with variations and special cases of the models defined in Section 2.

The remaining sections of the paper, excluding conclusion and references, deal with precise statements of our results and technical details of their proofs. For complex proofs, we first provide a more informal overview of the steps in the proof before proceeding with technical details. These sections are organized in the following manner.

- In Section 4 we discuss the quadratic integer programming reformulations of the various SBO problems.
- In Section 5 we state and prove our poly-logarithmic approximation algorithms for SSBO and MULTI-SSBO problems (main result **(R1)**).
- In Section 6, we state and prove our approximation-hardness results for both SSBO and MULTI-SSBO problems (main result **(R2)**).
- Section 7 contain all other results:
 - In Section 7.1 we show that many SSBO problems have improved solutions if certain parameters are restricted in their range of values.
 - In Section 7.2 we show the limitations of semidefinite programming based approaches for solving SSBO problems.
 -

2 Scenario Model for Stochastic Budget Optimization

We discuss the model and related problems using the language of sponsored search⁵. We use the suffix SSBO (Scenario Stochastic Budget Optimization) for various acronyms for different versions of our problems. For the convenience of the readers and to delay introducing more involved notations, we first start with a slightly simpler version of the model involving only one slot. We refer to this version as the “uniform cost” case and describe it in the next section.

2.1 Single Slot Case: Uniform Cost Model

This basic model starts with the following assumptions:

- There is a single slot for advertising.
- We have a set of n keywords $\mathcal{K}_1, \mathcal{K}_2, \dots, \mathcal{K}_n$ with the keyword \mathcal{K}_j having a *cost-per-click* d_j (a positive integer).
- We have a positive integer B denoting the *budget* for the advertiser.
- We have a collection of m “scenarios” where the i th scenario is characterized by the following parameters:
 - A probability of ε_i ($\sum_{i=1}^m \varepsilon_i = 1$).
 - A “click vector” $(a_{i,1}, a_{i,2}, \dots, a_{i,n})$ where each $a_{i,j} \geq 0$ is an integer. Each $a_{i,j}$ denotes the number of clicks obtained by the j th keyword \mathcal{K}_j in the i th scenario.

⁵Our discussion can easily be adopted to other internet ad channels like display ads and behavioral targeting.

Scenarios can be thought of as sampling the model over various times⁶.

Our general goal is to compute n selection variables x_1, x_2, \dots, x_n , where x_j corresponds to the j th keyword, to *maximize* a suitable total payoff. A crucial aspect of the discussed formulation is that, if the budget is not limiting, then the payoff corresponds to the total number of expected clicks, but if the budget turns out to be limiting for any scenario then the payoff *scales* the total number of expected clicks by the fraction that the budget would provide⁷. Based on the above intuition, our precise goal is *maximize* the total *expected* payoff over *all* scenarios, *i.e.*,

$$\text{maximize } \mathbb{E}[\text{payoff}] = \sum_{i=1}^m \mathbb{E}[\text{payoff}_i]$$

where the *expected payoff* $\mathbb{E}[\text{payoff}_i]$ for the i th scenario is

$$\mathbb{E}[\text{payoff}_i] = \begin{cases} \varepsilon_i \sum_{j=1}^n a_{i,j} x_j, & \text{if } \sum_{j=1}^n a_{i,j} d_j x_j \leq B \\ \frac{B}{\sum_{j=1}^n a_{i,j} d_j x_j} \left(\varepsilon_i \sum_{j=1}^n a_{i,j} x_j \right), & \text{otherwise} \end{cases} \quad (1)$$

Following [21], we distinguish between two versions of the problem based on the nature of the selection variables:

Integral version (UNIFORM-INT-SSBO): $x_j \in \{0, 1\}$ for all j . This corresponds to the case when based on the stochastic information, either the advertiser chooses to win and pay for all clicks for a keyword, or not at all. Hence, the strategy of the advertiser is *deterministic*.

Fractional version (UNIFORM-FRAC-SSBO): $0 \leq x_j \leq 1$ for all j . This can be thought of as a strategy in which the advertiser treats these numbers as *probabilities* and bids for the keywords in a *randomized* fashion based on these probabilities, thereby only winning (and paying for) a portion of all clicks and impressions for each keyword. If the deterministic strategy is hard to compute and provides a solution of bad quality then the randomized strategy is more desirable.

Other than the scenario model, there are at least two other possible models for stochastic budget optimization as discussed in [21]. In the *proportional model* there is just one global random variable for the total number of clicks in the day that keeps the relative proportions of clicks for different keywords the same, whereas in the *independent keywords model* each keyword comes with its own probability distribution. However, among all these models this scenario-based model is perhaps one of the most natural model of reality and provides an appropriate middle ground between complex arbitrary joint probability distribution and a single distribution for all keywords. It was shown in [21] that both UNIFORM-INT-SSBO and UNIFORM-FRAC-SSBO are NP-hard. In the sequel, we assume without loss of generality that $1 = d_1 \leq d_2 \leq \dots \leq d_n$.

⁶Scenarios can be provided by the search engine for the advertisers, or used by the search engines to bid on behalf of advertisers. Similarly, advertisers and other search engine optimizers can also “infer” scenarios indirectly using trends and other data provided by search engines.

⁷The underlying assumption is that, within a scenario, the queries and keywords are well-mixed and, when budget runs out, the ad campaign is halted for the period as is currently done. The queries and keywords are well-mixed not only because of aggregation of streams from millions of users but also because of ad throttling that spreads out the eligible ad campaigns over the period of a scenario. See [21] for exact details of justification.

2.2 Single Slot Case: General Model

In a more realistic version of the SSBO problems the cost-per-click values may vary *slightly* over a range of scenarios due to their small errors in estimation. This can be modeled by introducing a stretch parameter (*small integer*)⁸ $1 \leq \kappa = O(\text{poly}(\log(m+n)))$. Now, d_j stands for the *basic* cost-per-click for the keyword \mathcal{K}_j , whereas the *real* cost-per-click for the keyword \mathcal{K}_j in the i th scenario is denoted by $c_{i,j}$, with $c_{i,j} \in [d_j, \kappa d_j]$ ⁹. Then, Equation (1) can be simply updated by replacing d_j in the equation of the i th scenario by $c_{i,j}$. We refer to the integral and fractional versions of this general case as INT-SSBO and FRAC-SSBO, respectively; note that the UNIFORM-SSBO problems are obtained from the corresponding SSBO problems by setting $\kappa = 1$.

2.3 Multi Slot Model

In the multi-slot case there are $s \geq 1$ *slots* for each keyword with the *GSP second price auction* for these slots. Let $d_{j,k}$ be an integer denoting the value of the *basic* cost-per-click the k th slot of the j th keyword; we assume without loss of generality that $d_{j,1} \leq d_{j,2} \leq \dots \leq d_{j,s}$. Let $c_{i,j,k} \in [d_{j,k}, \kappa d_{j,k}]$ denote the value of the *real* cost-per-click for the k th slot of the j th keyword in the i th scenario where κ is the stretch parameter as in Section 2.2, and let $B > 0$ denote the *budget* (a positive integer) for the advertiser. Our goal is now to compute a set of sn selection variables $x_{j,k}$ where the selection variable $x_{j,k}$ corresponds to k th slot for the j th keyword. We again have a collection of m scenarios where the i th scenario is characterized via:

- a probability ε_i ($\sum_{i=1}^m \varepsilon_i = 1$), and
- a “click vector” $(a_{i,j,1}, a_{i,j,2}, \dots, a_{i,j,s})$ where each $a_{i,j,k}$ is a non-negative integer;

The goal is to compute the allocation variables $x_{j,k}$ ’s with the constraints

$$\forall j : \sum_{k=1}^s x_{j,k} \leq 1 \quad (2)$$

to maximize the total expected payoff

$$\mathbb{E}[\text{payoff}] = \sum_{i=1}^m \mathbb{E}[\text{payoff}_i]$$

where

$$\mathbb{E}[\text{payoff}_i] = \begin{cases} \varepsilon_i \sum_j \sum_k a_{i,j,k} x_{j,k}, & \text{if } \sum_j \sum_k a_{i,j,k} c_{i,j,k} x_{j,k} \leq B \\ \frac{B}{\sum_j \sum_k a_{i,j,k} c_{i,j,k} x_{j,k}} \left(\varepsilon_i \sum_j \sum_k a_{i,j,k} x_{j,k} \right), & \text{otherwise} \end{cases} \quad (3)$$

We again distinguish between two versions of the problem:

Integral version (INT-MULTI-SSBO): $x_{j,k} \in \{0, 1\}$ for all j and k . Here, $x_{j,k} = 1$ if the advertiser selects the k th slot for the j th keyword, and $x_{j,k} = 0$ otherwise.

⁸Throughout the paper, the notation $\text{poly}(a)$ denotes a polynomial in a , *i.e.*, a^c for some positive constant c .

⁹For example, the stretch parameter κ allows us to model situations such as when the real costs can be drawn from a probability distribution with a mean around $\frac{1+\kappa}{2}d_j$ with a negligible probability of occurring outside a range of $\pm \frac{1-\kappa}{2}d_j$ of the mean. Note that this is just an illustration. We do not assume any specific probability distribution for the variations of the real costs per click except that it varies within an interval of length κ .

Fractional version (FRAC-MULTI-SSBO): $0 \leq x_{j,k} \leq 1$ for all j and k . Here, $x_{j,k}$ denotes the probability that the advertiser selects the k th slot for the j th keyword and $1 - (\sum_{k=1}^s x_{j,k})$ is the probability with which the advertiser does not bid on the j th keyword at all.

Note that the scenario model for multi-slot stochastic budget optimization is *quite different* in nature from the other multi-slot models such as the one discussed in [12] since, for example, one can go under or over the budget in one scenario to get a better overall expected payoff.

2.4 Relevance and Significance of Scenario Models

Scenario models are a popular way of modeling optimization problems involving uncertainties in parameters by creating a number of scenarios that depict the probability distribution of various possibilities and then provide a solution that optimizes the expectations of outcomes over these scenarios. The scenario model is important for at least two reasons as explained in [21], which we state below. Firstly, market analysts often think of uncertainty by *explicitly* creating a set of a few model scenarios, possibly attaching a weight to each scenario. Secondly, the scenario model gives us an important tool into understanding the fully general problem with arbitrary joint distributions. Allowing the full generality of an arbitrary joint distribution gives us significant modeling power, but poses challenges to the algorithm designer. Since a naive explicit representation of the joint distribution requires space exponential in the number of random variables, one often represents the distribution implicitly by a sampling oracle. A common technique, Sampled Average Approximation, is to replace the true distribution by a uniform or non-uniform distribution over a set of samples drawn by some process from the sampling oracle, effectively reducing the problem to the scenario model. In addition to their usual applications in operations research (*e.g.*, see [8]), this approach is getting more and more attention in Wall Street as financial portfolios are being created in this way (*e.g.*, see [25]). For example, Cocco, Consiglio and Zenios in [7] developed a scenario-based optimization model for asset and liability management of participating insurance policies with minimum guarantees and Mausser and Rosen in [16] developed three scenario optimization models for portfolio credit risk.

In sponsored search, this is an appropriate model and embodies the “best response” strategy. There is a complex function that maps the state of the world and the users to the queries they pose and their actions such as whether they click on ads. The search engines give a limited amount of information to help advertisers¹⁰, and advertisers can learn various scenarios that determine their click vs cost behaviors to some extent by running experiments, analyzing their web traffic *etc.* However, sponsored search products only provide a limited bidding language to structure one’s campaign¹¹ and hence, necessarily, most advertisers have to target different scenarios simultaneously with each bidding choice. This is the stochastic budget optimization problem we study in this paper. One natural idea is for advertisers to recognize in real time the particular scenario one faces and then apply the best bidding for that scenario. However, this is difficult to do in practice because of limited and delayed information in the system, and it is also expensive to implement. Thus, stochastic budget optimization problems under the scenario model are very appropriate for sponsored search applications.

We do acknowledge that other strategies besides the “best response” may be used by advertisers in practice¹², and stochastic budget optimization algorithms proposed here are not currently used within the practical tools that are publicly available. Nevertheless, best response is a reasonable

¹⁰For example, see <https://adwords.google.com/select/TrafficEstimatorSandbox>

¹¹See for example, <http://algo.research.googlepages.com/ec09-partI.pdf>

¹²By other strategies, we mean strategies in which the advertiser does not fix the strategies of other advertisers.

strategy (even recommended by some search engines), and indeed many anecdotal conversations with advertisers and sponsored search optimizers have clearly indicated to us that they would like to bid to balance across myriad of scenarios. Our algorithms in this paper (even the dynamic programming based ones) can be *easily* implemented in current systems.

2.5 Notational Remarks

As the reader may have already observed, precise definitions of the various models involve a lot of variables and subscripts. To make the exposition clearer, we will therefore adopt the following conventions:

- For variables involving keywords, scenarios and (for the multi-slot model) slots, we will use subscripts i, j and k (and their obvious variations such as $i_1, i', \text{etc.}$) for scenarios, keywords and slots, respectively.
- Variables such as $m, n, \mathcal{K}_j, d_j, B, \varepsilon_i, a_{i,j}, a_{i,j,k}, c_{i,j}, c_{i,j,k}, x_j, x_{j,k}, \text{payoff}, \text{payoff}_i, \kappa, s$ and B , when used in the context of the stochastic budget optimization models, will be used for their intended meanings as described in Sections 2.1–2.3.
- Note that:
 - $m, n, d_j, B, a_{i,j}, a_{i,j,k}, c_{i,j}, c_{i,j,k}$ and s are positive integers;
 - $0 \leq \varepsilon_i \leq 1$ and $\sum_{i=1}^m \varepsilon_i \leq 1$;
 - $1 \leq \kappa = O(\text{poly}(\log(m+n)))$ is an integer. We refer to this in the sequel by the phrase “ κ is a small integer”.
- The size of an input instance of our SBO problems, which we will denote by **size-of-input** and which is *crucial in differentiating polynomial-time algorithms from pseudo-polynomial-time algorithms*, is as follows:

- For INT-SSBO and FRAC-SSBO:

$$\text{size-of-input} = \text{poly} \left(m + n + \left(\max_{\substack{1 \leq i \leq m \\ 1 \leq j \leq n}} \log_2 a_{i,j} \right) + \left(\max_{\substack{1 \leq i \leq m \\ 1 \leq j \leq n}} \log_2 c_{i,j} \right) + \left(\max_{1 \leq i \leq m} \frac{1}{\varepsilon_i} \right) \right).$$

- For INT-MULTI-SSBO and FRAC-MULTI-SSBO,

$$\text{size-of-input} = \text{poly} \left(s + m + n + \left(\max_{\substack{1 \leq i \leq m \\ 1 \leq j \leq n \\ 1 \leq k \leq s}} \log_2 a_{i,j,k} \right) + \left(\max_{\substack{1 \leq i \leq m \\ 1 \leq j \leq n \\ 1 \leq k \leq s}} \log_2 c_{i,j,k} \right) + \left(\max_{1 \leq i \leq m} \frac{1}{\varepsilon_i} \right) \right).$$

On rare occasions, if we need to reuse the above-mentioned indices or variables and thus deviate from these conventions, the accompanying text will make the deviation clear.

3 Summary of Results and Proof Techniques

[21] left the computational complexity issues of the scenario model as the main open problem after showing that both the integral and fractional versions of this problem, even for single slot case,

are NP-hard and noting that *no* non-trivial approximability results are known. While prior results for (S)BO problems exploit insights from the Knapsack problem to associate some potential payoff with each keyword, a central difficulty encountered in directly applying those techniques for our models is that payoff from a keyword can be very different from one scenario to another.

3.1 Summary of Results

We provide a slightly coarse summary of the results obtained in this paper; precise bounds are available in the corresponding technical section that proves the result.

Main Results

(R1) (Approximation algorithms): We provide algorithms that run in *near-linear* time and achieve the following approximation ratios¹³:

- $\min \{O(m), O(\kappa \log d_n)\}$ -approximation for both INT-SSBO and FRAC-SSBO and,
- $\min \{O(m), O(s \kappa \log \Delta \log^2(m+n))\}$ -approximation for INT-MULTI-SSBO and FRAC-MULTI-SSBO, where $\Delta = \max_{j,k} d_{j,k}$.

(R2) (Approximation hardness for the single slot case) We show that, unless $ZPP = NP$, there exists instances of INT-SSBO and FRAC-SSBO, with n keywords and $m = n$ scenarios each with equal probability, such that *any* polynomial-time algorithm for solving these problems must have an approximation ratio of any one of the following (for any constant $0 < \varepsilon < 1$):

- $\Omega(m^{1-\varepsilon})$ (and, thus, also $\Omega(n^{1-\varepsilon})$), or
- $\Omega(\kappa \log^{1-\varepsilon} d_n)$.

This *almost matches* the upper bounds in **(R1)**. Thus, we cannot in general improve the approximation bound in **(R1)**.

(R3) (Approximation hardness for the multi-slot case) Since SSBO problems are special case of MULTI-SSBO problems for $s = 1$, the approximation hardness bounds for SSBO can be extended to MULTI-SSBO in an obvious manner. However, because of the fact that selection variables of various slots of the same keyword are dependent on each other via constraints such as Equation (2), our lower bound proofs translate to corresponding lower bounds of the form $\Omega(m^{1-\varepsilon})$, $\Omega(n^{1-\varepsilon})$, or $\Omega(\log \kappa \cdot \log^{1-\varepsilon} d_n)$ for an MULTI-SSBO instance with n keywords, $m = n$ scenarios and s slot. Thus, unfortunately, the lower bounds are *independent* of s , though one would expect the computational complexity of the problem to depend on s , say when s is large compared to n .

Using a different amplification of NP-hard problems as suggested by Raz's parallel repetition theorem [11, 18, 22], we can show that there exists instances of INT-MULTI-SSBO and FRAC-MULTI-SSBO with n keywords and $s > 1$ slots such that *any* polynomial-time algorithm for solving these problems must have an approximation ratio of $2^{\log^{1-\varepsilon}(ns)}$ for any constant $0 < \varepsilon < 1$, provided $NP \not\subseteq DTIME(n^{\text{poly}(\log n)})$ ¹⁴. As an example, $2^{\log^{1-\varepsilon}(ns)}$ dominates $n^{1-\varepsilon}$ if $s = \Omega(n^{\log n})$.

¹³The reader is reminded that $\kappa = O(\text{poly}(\log(m+n)))$.

¹⁴More detailed discussions on the generalization of the lower bound for the single-slot case to the multi-slot case and comparison of the two lower bounds is available in Section 6.2.

We also show that INT-MULTI-SSBO is MAX-SNP-hard for $s = 2$ even when $\kappa = 1$ and $c_{j,k} = 1$ for all j and k .

Other Results

In addition to the main results, we also prove a number of other results dealing with variations and special cases of our problems.

Fixed parameter tractability issues: For certain parameter ranges of practical interest we show that these optimizations problems can be solved efficiently. If m or ns is fixed, FRAC-MULTI-SSBO has a polynomial time solution with an absolute error of δ for any fixed $\delta > 0$. If additionally bids are polynomial in size, INT-MULTI-SSBO also has a polynomial time solution with an absolute error of δ for any fixed $\delta > 0$.

Limitations of semi-definite programming based approaches: The lower bounds in (R2) have $\varepsilon < 1$ and thus leaves a “very small” gap between this lower bound and the upper bounds described in (R1). It is natural to ask if the gap could be eliminated; for example can we design an approximation algorithm for the special case for $\kappa = 1$ whose approximation ratio is, say, $o\left(\frac{m}{\log m}\right)$ or $o\left(\frac{\log d_n}{\log \log d_n}\right)$? Although we are unable to provide a concrete proof that such a polynomial time approximation algorithm does not exist, we nonetheless observe that the natural semidefinite programming relaxation will not work since it has a large integrality gap of $\frac{m}{2} = \Theta\left(\frac{\log d_n}{\log \log d_n}\right)$.

Dual of SSBO problems: Finally, in some cases, the *dual* of the stochastic budget optimization problem may be of interest, where we are given a target expected number of clicks and the goal is to minimize the expected budget spent while reaching the target. We present some exact and approximate results for this dual version of the problem.

3.2 Brief Overview of Proof Techniques

In general, budget optimization problems are akin to knapsack problems¹⁵. But the stochastic budget optimization problems studied in this paper are different because their budgets are “soft”, *i.e.*, they can be exceeded, if under a suitable scaling they meet the budget constraint, and this improves the objective function. The stochastic budget optimization problems can be more insightfully thought of as special bipartite quadratic programs (these with ± 1 variables correspond to Grothendieck’s inequality with a nice history, but we have 0/1 variables). Standard approaches to solving other special cases of quadratic programs, for example, using relaxations via semi-definite programming, do not provably work as we show. Instead, for upper bounds, we take alternative combinatorial approaches. For showing hardness results, we use intuitions from connections of our problems to these quadratic programs. For one proof, we show reduction from the hard instances of the *maximum independent set* problem [15] on graphs to the bipartite 0/1 quadratic integer programming reformulations of FRAC-SSBO and INT-SSBO. For the hardness proofs for multi-slot case, we have to start from an inapproximability result of certain type of multi-prover systems to obtain the best hardness results, again crucially using the connection to these quadratic programs. While anecdotally one may indeed believe these problems to be computationally hard, our results show that this is not true for many ranges of parameters of interest, but do identify the parameter settings that make them computationally hard. Taken together, our results are the *first known* non-trivial

¹⁵See for example, http://algo.research.googlepages.com/ec09_pub.pdf

complexity results for stochastic budget optimization problems under the scenario model beyond NP-hardness.

4 SBO Problems and Bipartite Quadratic Integer Programs

In this section we show how to reformulate various SBO problems as *bipartite quadratic integer programs* (QIP). These reformulations are heavily used in later proofs in the paper. A bipartite quadratic program is a quadratic program in which there is a bipartition of variables such that every term involves at most one variable from each partition. A well-known example of such a (strict) quadratic program on variables taking ± 1 values is the so-called Grothendieck's inequality [1]. However, as will show later, our quadratic program differs significantly in nature from this inequality.

4.1 SSBO and QIP

<p>(* Quadratic program (Q1) *)</p> <p>(* $w_{i,j} = y_{i,j} c_{i,j}$ for all i and j *)</p> <p>maximize $\sum_{i=1}^m \sum_{j=1}^n \alpha_i x_j y_{i,j}$</p> <p>subject to</p> $\forall 1 \leq i \leq m: \alpha_i \left(\sum_{j=1}^n w_{i,j} x_j \right) \leq B_i$ $\forall 1 \leq i \leq m: 0 \leq \alpha_i \leq 1$ $\forall 1 \leq j \leq n: 0 \leq x_j \leq 1$	<p>(* Quadratic program (Q2) *)</p> <p>(* $w_{i,j,k} = y_{i,j,k} c_{i,j,k}$ for all i, j and k *)</p> <p>maximize $\sum_{i=1}^m \alpha_i \left(\sum_{j=1}^n \sum_{k=1}^s x_{j,k} y_{i,j,k} \right)$</p> <p>subject to</p> $\forall 1 \leq i \leq m: \alpha_i \left(\sum_{j=1}^n \sum_{k=1}^s w_{i,j,k} x_{j,k} \right) \leq B_i$ $\forall 1 \leq j \leq n: \sum_{k=1}^s x_{j,k} \leq 1$ $\forall 1 \leq i \leq m: 0 \leq \alpha_i \leq 1$ $\forall 1 \leq j \leq n \forall 1 \leq k \leq s: 0 \leq x_{j,k} \leq 1$
---	--

Figure 1: Quadratic Integer Programs for SBO Problems. Y is a matrix with non-negative entries ($y_{i,j}$ for **(Q1)** and $y_{i,j,k}$ for **(Q2)**) and B_1, B_2, \dots, B_m are positive real numbers.

We show how to reformulate SSBO as a bipartite quadratic integer program. Consider the quadratic program **(Q1)** in Fig. 1. By “integral version” of **(Q1)** we refer to replacing the constraint $0 \leq x_i \leq 1$ by $x_i \in \{0, 1\}$.

Proposition 1. *The quadratic program (Q1) and its integral version is equivalent to INT-SSBO or FRAC-SSBO, respectively.*

Proof. Consider an instance of SSBO. Let $y_{i,j} = \varepsilon_i a_{i,j}$, $w_{i,j} = c_{i,j} y_{i,j}$ and $B_i = \varepsilon_i B$. Intuitively, the B_i 's correspond to budgets for the i th scenario scaled by the probability of the i th scenario, and the $y_{i,j}$'s are the $c_{i,j}$'s scaled by the probability of the i th scenario. Then, the inequality $\sum_{j=1}^n a_{i,j} c_{i,j} x_j \leq B$ becomes $\sum_{j=1}^n y_{i,j} c_{i,j} x_j \leq B_i \equiv \sum_{j=1}^n w_{i,j} x_j \leq B_i$ and the fraction $\frac{B}{\sum_{j=1}^n a_{i,j} c_{i,j} x_j}$ becomes $\frac{B_i}{\sum_{j=1}^n w_{i,j} x_j}$. Conversely, given an instance of **(Q1)**, let $B = \sum_{i=1}^m B_i$, $\varepsilon_i = \frac{B_i}{B}$ and $a_{i,j} = \frac{y_{i,j}}{\varepsilon_i}$. Thus, $\varepsilon_i a_{i,j} = y_{i,j}$, the inequality $\sum_{j=1}^n w_{i,j} x_j \leq B_i \equiv \sum_{j=1}^n y_{i,j} c_{i,j} x_j \leq B_i$ is the same as $\sum_{j=1}^n a_{i,j} c_{i,j} x_j \leq B$ and the fraction $\frac{B_i}{\sum_{j=1}^n w_{i,j} x_j}$ is the same as $\frac{B}{\sum_{j=1}^n a_{i,j} c_{i,j} x_j}$. Thus, in the sequel, we assume such a correspondence.

Now, consider a solution vector $\mathbf{x} = (x_1, x_2, \dots, x_n)$ and $\boldsymbol{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_m)$ for **(Q1)**. Then \mathbf{x} also defines a solution vector for SSBO. We must verify that this is indeed a valid solution vector with a correct expected payoff. Let $Q_i = \sum_{j=1}^n w_{i,j} x_j$. If $\alpha_i Q_i < B_i$ then $\alpha_i = 1$ since otherwise the solution for **(Q1)** can be further improved, and then $\mathbb{E}[\text{payoff}_i] = \sum_{j=1}^n y_{i,j} x_j$, which is correct. Otherwise $\alpha_i Q_i = B_i$ and then $\mathbb{E}[\text{payoff}_i] = \frac{B_i}{B_i/\alpha_i} \sum_{j=1}^n y_{i,j} x_j = \alpha_i \sum_{j=1}^n y_{i,j} x_j$, which is also correct. This shows that for every instance of **(Q1)** there is a corresponding instance of SSBO with the same expected payoff.

Now, consider a solution vector \mathbf{x} for SSBO. Then, if $Q_i \geq B_i$ then $\alpha_i = B_i/Q_i$ otherwise $\alpha_i = 1$. It is easy to see in the same manner that this provides a valid solution of **(Q1)** with the same objective value. \square

Relationship to the Standard Knapsack Problems

If $m = \kappa = 1$ and α_1 is set to a fixed constant, then **(Q1)** reduces a special linear program which is equivalent to the so-called (fractional) knapsack problem which is well-studied in the literature. Extending this analogy, by the phrase “the standard fractional knapsack problem corresponding to the i th row of and p th through q th column of Y ”, we will mean the linear program as shown in Fig. 2 (it is easy to see that there is an optimal solution of this

linear program in which $\alpha_i = \min \left\{ 1, \frac{B_i}{\left(\sum_{j=p}^q c_{i,j} y_{i,j} \right)} \right\}$). Since $d_p \leq d_{p+1} \leq \dots \leq d_q$, the following well-known fact follows.

Fact 1. [13] *An optimal solution to the linear program in Fig. 2 (“optimal payoff for the i th row and p th through q th column of Y ”) is a “prefix solution”, i.e., there is an index j' such that $x_j = 1$ for $j < j'$, $0 < x_{j'} \leq 1$ and $x_j = 0$ for $j > j'$.*

$$\begin{array}{ll} \text{maximize} & \alpha_i \sum_{j=p}^q y_{i,j} x_j \\ \text{subject to} & \alpha_i \left(\sum_{j=p}^q c_{i,j} y_{i,j} x_j \right) \leq B_i \\ & \forall p \leq j \leq q: 0 \leq x_j \leq 1 \end{array}$$

Figure 2: LP for i th row and p th through q th column of Y .

4.2 MULTI-SSBO and QIP

The quadratic programming reformulation of MULTI-SSBO can also be obtained in a similar manner and is shown as **(Q2)** in Fig. 1.

5 Poly-logarithmic Approximations for SSBO and MULTI-SSBO (main result (R1))

Theorem 1 (Near-linear time approximation). *There is a*

- (i) $\min \{O(m), O(\kappa \log d_n)\}$ -approximation for both INT-SSBO and FRAC-SSBO;
- (ii) $\min \{O(m), O(s\kappa \log \Delta)\}$ -approximation for FRAC-MULTI-SSBO and
- (iii) $\min \{O(m), O(s\kappa \log \Delta \log^2(m+n))\}$ -approximation for INT-MULTI-SSBO

where, for (ii) and (iii), $\Delta = \max_{j,k} d_{j,k}$. All these algorithms can be implemented in linear or near-linear time using standard data structures and algorithmic techniques.

-
1. Partition the keywords into *maximal* groups such that if a group G contains p th through q th keyword then $d_q/d_p \leq 2$ and $d_{q+1}/d_p > 2$.
Let \mathcal{G} be the set of such groups.
 2. **For** each group $G \in \mathcal{G}$ consisting of keywords, say $\mathcal{K}_p, \mathcal{K}_{p+1}, \dots, \mathcal{K}_q$, **do**
 Set $x_j = 1$ for every $p \leq j \leq q$;
 let $\mathbb{E}[\text{payoff}']$ be the payoff of this solution
 3. Output the best of the solutions obtained in 2.
-

Figure 3: Algorithm for the case of $\kappa = 1$.

In the rest of this section, we prove the above theorem. As a first attempt, one might be tempted to use recent techniques in designing efficient algorithms for multiple-knapsack problems [6, 17] for our problem; however it is not difficult to design examples where such approaches fail badly since our budget constraints are “soft” (they can be exceeded if scaling them gives better payoff) and our probabilities are “arbitrary”. As a second attempt, one might take our quadratic programming reformulation as discussed in Section 4 and semidefinite-programming based rounding approach such as in [14]. However, it can be shown that the integrality gap of such a reformulation is very large. The failure of these natural approaches shows the difficulty of the problems. Thus, we are led to explore other combinatorial approaches to provide the desired approximation.

5.1 $O(m)$ -approximation for INT-SSBO and FRAC-SSBO

To get a $O(m)$ -approximation we can do the following. For each i we solve the standard (integer or fractional) knapsack problem for the i th row of Y ; let p_i be the value of an optimal solution. Then, take the best of these solutions, say of value $p = \max_{1 \leq i \leq m} \{p_i\}$. Each fractional knapsack problem can be solved exactly in $O(n \log n)$ time [13] and a $O(n \log n)$ time greedy 2-approximation algorithm for the integer knapsack problem is also well known [19].

We now note that $\mathbb{E}[\text{payoff}] \leq \sum_{i=1}^m p_i$. Indeed, consider an optimal solution of SSBO. If $\alpha_i = 1$, then by definition of p_i we have $\mathbb{E}[\text{payoff}_i] \leq p_i$. If $\alpha_i < 1$, then we set $\alpha_i = 1$ and set a new value of x_j as $x'_j = \alpha_i x_j$. This does not change $\mathbb{E}[\text{payoff}_i]$ and now we again have $\mathbb{E}[\text{payoff}_i] \leq p_i$. Thus, we have $p \geq \mathbb{E}[\text{payoff}]/m$.

If $p = p_i$ for some i , then the solution of the knapsack problem of value p can be extended to a solution of SSBO by setting $\alpha_{i'} = 0$ for $i' \neq i$.

5.2 $O(\kappa \log d_n)$ -approximation for INT-SSBO and FRAC-SSBO

Case of $\kappa = 1$: Uniform Cost Model

The algorithm is shown in Fig. 3. Consider a group $G \in \mathcal{G}$ consisting of the keywords $\mathcal{K}_p, \mathcal{K}_{p+1}, \dots, \mathcal{K}_q$. By the “SSBO problem on G ” we mean the instance of the SSBO problem in which our click in-

put consists of the submatrix $Y_{p,q} = \begin{pmatrix} y_{1,p} & y_{1,p+1} & \cdots & y_{1,q} \\ \vdots & \vdots & \cdots & \vdots \\ y_{m,p} & y_{m,p+1} & \cdots & y_{m,q} \end{pmatrix}$ of Y containing all rows and p th

through q th columns, the costs-per-click d_p, \dots, d_q , the budgets B_1, \dots, B_m , and the selection variables x_p, \dots, x_q . Let $\mathbb{E}[\text{payoff}_G]$ be the value of expected payoff of an optimal solution for this subproblem. Since $\max_{G \in \mathcal{G}} \mathbb{E}[\text{payoff}_G] \geq \frac{\mathbb{E}[\text{payoff}]}{|\mathcal{G}|}$ and $|\mathcal{G}| = O(\log d_n)$, the following lemma proves the desired approximation bound.

Lemma 2. $\mathbb{E}[\text{payoff}'] \geq \frac{\mathbb{E}[\text{payoff}_G]}{2}$.

Proof. We only need to prove the lemma for the case when $\mathbb{E}[\text{payoff}_G]$ is the total expected payoff of an optimal solution of the FRAC-SSBO problem on G since obviously the total expected payoff of an optimal solution of the INT-SSBO problem on G is no more than $\mathbb{E}[\text{payoff}_G]$. Let $D = \sum_{j=p}^q d_j y_{i,j}$ and $\beta = |G|$. By our choice of the group G ,

$$\beta d_p \sum_{j=p}^q y_{i,j} \leq D \leq \beta d_q \sum_{j=p}^q y_{i,j} \leq 2\beta d_p \sum_{j=p}^q y_{i,j}.$$

Using the quadratic programming formulation **(Q1)** and remembering that $c_{i,j} = d_j$ when $\kappa = 1$, the FRAC-SSBO instance on G is equivalent to the following quadratic program **(Q3)**:

(* Quadratic program (Q3) *)

$$\begin{aligned} & \text{maximize } \sum_{i=1}^m \alpha_i \left(\sum_{j=p}^q y_{i,j} x_j \right) \\ & \text{subject to } \forall 1 \leq i \leq m: \alpha_i \left(\sum_{j=p}^q d_j y_{i,j} x_j \right) \leq B_i \\ & \quad \forall 1 \leq i \leq m: 0 \leq \alpha_i \leq 1 \\ & \quad \forall p \leq j \leq q: 0 \leq x_j \leq 1 \end{aligned}$$

Fix any optimal solution for our FRAC-SSBO instance on G , *i.e.*, fix an optimal solution vector $(\alpha_1^*, \alpha_2^*, \dots, \alpha_m^*)$ and $(x_p^*, x_{p+1}^*, \dots, x_q^*)$ of **(Q3)**. In our solution sets $x_p = x_{p+1} = \dots = x_q = 1$; thus $\alpha_i = \min \left\{ 1, \frac{B_i}{D} \right\}$ for every i and $x_j \geq x_j^*$ for every $p \leq j \leq q$.

Case 1: $D \leq B_i$. Then, $\alpha_i = 1 \geq \alpha_i^*$, $x_j = 1 \geq x_j^*$ for $p \leq j \leq q$, and thus

$$\alpha_i \left(\sum_{j=p}^q y_{i,j} x_j \right) \geq \alpha_i^* \left(\sum_{j=p}^q y_{i,j} x_j^* \right).$$

Case 2: $D > B_i$. Then, $\alpha_i = \frac{B_i}{D}$. Let α_i^+ and $(x_p^+, x_{p+1}^+, \dots, x_q^+)$ be the optimal prefix solution of the i th row and p th through q th column of Y . By Fact 1, there exists an index $1 \leq \ell \leq \beta$ such that $x_j^+ = 1$ for $j < p + \ell - 1$, $0 < x_{p+\ell-1}^+ \leq 1$ and $x_j^+ = 0$ for $j > p + \ell - 1$, and

$$B_i = \alpha_i^+ \sum_{j=p}^{p+\ell-1} d_j y_{i,j} x_j^+ = \alpha_i^+ \sum_{j=p}^q d_j y_{i,j} x_j^+ \geq \alpha_i^* \sum_{j=p}^q d_j y_{i,j} x_j^*$$

Now, it follows that

$$\begin{aligned} \alpha_i \left(\sum_{j=p}^q y_{i,j} x_j \right) &= \left(\frac{B_i}{D} \right) \times \sum_{j=p}^q y_{i,j} \geq \left(\frac{B_i}{D} \right) \times \left(\frac{\sum_{j=p}^q d_j y_{i,j}}{d_q} \right) = \frac{B_i}{d_q} = \frac{\alpha_i^+ \sum_{j=p}^{p+\ell-1} d_j y_{i,j} x_j^+}{d_q} \\ &\geq \frac{d_p}{d_q} \times \left(\alpha_i^+ \sum_{j=p}^{p+\ell-1} y_{i,j} x_j^+ \right) \geq \frac{1}{2} \times \left(\alpha_i^+ \sum_{j=p}^{p+\ell-1} y_{i,j} x_j^+ \right) \geq \frac{1}{2} \times \alpha_i^* \left(\sum_{j=p}^q y_{i,j} x_j^* \right) \end{aligned}$$

Thus, combining both cases, we have

$$\mathbb{E}[\text{payoff}'] = \sum_{i=1}^m \alpha_i \left(\sum_{j=p}^q y_{i,j} x_j \right) \geq \frac{1}{2} \times \sum_{i=1}^m \alpha_i^* \left(\sum_{j=p}^q y_{i,j} x_j^* \right) = \frac{\mathbb{E}[\text{payoff}_G]}{2}$$

□

Case of $\kappa > 1$: General Single-slot Model

Using our δ -approximation algorithm for UNIFORM-SSBO (for $\delta = O(\log d_n)$) as outlined in Fig. 3, we show how to use it as a subroutine to get a $\kappa \delta = O(\kappa \log d_n)$ -approximation for INT-SSBO (and, hence, also for FRAC-SSBO). The algorithm is shown in Fig. 4.

-
1. Replace (truncate) each $c_{i,j}$ by its *new* value $c'_{i,j} = d_j$.
 2. Use the approximation algorithm in Fig. 3 with these new truncated values of $c_{i,j}$'s. Let $\mathbf{x} = (x_1, x_2, \dots, x_n)$ and $\boldsymbol{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_m)$ be the solution vectors returned.
 3. Output \mathbf{x} and $\boldsymbol{\alpha}' = (\alpha'_1, \alpha'_2, \dots, \alpha'_m) = (\frac{\alpha_1}{\kappa}, \frac{\alpha_2}{\kappa}, \dots, \frac{\alpha_m}{\kappa})$ as our solution.
-

Figure 4: $O(\kappa \log d_n)$ -approximation algorithm for INT-SSBO.

We use the following notations:

- $\mathbf{x}^* = (x_1^*, x_2^*, \dots, x_n^*)$ and $\boldsymbol{\alpha}^* = (\alpha_1^*, \alpha_2^*, \dots, \alpha_m^*)$ be the solution vectors for an optimal solution of our (original) instance of SSBO, and $\mathbb{E}[\text{payoff}^*] = \sum_{i=1}^m \alpha_i^* \left(\sum_{j=1}^n y_{i,j} x_j^* \right)$ is the total expected payoff of this optimal solution.
- $\mathbf{x}^+ = (x_1^+, x_2^+, \dots, x_n^+)$ and $\boldsymbol{\alpha}^+ = (\alpha_1^+, \alpha_2^+, \dots, \alpha_m^+)$ be the solution vectors for an optimal solution of the truncated instance of SSBO, and $\mathbb{E}[\text{payoff}^+] = \sum_{i=1}^m \alpha_i^+ \left(\sum_{j=1}^n y_{i,j} x_j^+ \right)$ is the total expected payoff of this optimal solution.
- $\mathbb{E}[\text{payoff}] = \sum_{i=1}^m \alpha'_i \left(\sum_{j=1}^n y_{i,j} x_j \right)$ is the total expected payoff of the solution obtained by using the algorithm in Fig. 4.

Proposition 2. *The following statements are true:*

- (a) \mathbf{x} and $\boldsymbol{\alpha}'$ correspond to a valid solution of the SSBO instance.
- (b) $\mathbb{E}[\text{payoff}^+] \geq \mathbb{E}[\text{payoff}^*]$.
- (c) $\mathbb{E}[\text{payoff}] \geq \frac{\mathbb{E}[\text{payoff}^+]}{\kappa}$.

Thus the algorithm in Fig. 4 is a $O(\kappa \log d_n)$ -approximation.

Proof.

- (a) $\alpha'_i c_{i,j} = \frac{\alpha_i}{\kappa} c_{i,j} \leq \frac{\alpha_i}{\kappa} \kappa d_j = \alpha_i d_j$, thus $\alpha_i \left(\sum_{j=1}^n y_{i,j} d_j x_j \right) \leq B_i$ implies $\alpha'_i \left(\sum_{j=1}^n y_{i,j} c_{i,j} x_j \right) \leq B_i$.
- (b) The solution vectors \mathbf{x}^* and $\boldsymbol{\alpha}^*$ for an optimal solution of the SSBO instance is also a valid (not necessarily optimal) solution vector for the truncated instance of SSBO since $c'_{i,j} \leq c_{i,j}$.
- (c) This follows since $\alpha'_i = \frac{\alpha_i}{\kappa}$. □

5.3 Approximation Bounds for FRAC-MULTI-SSBO and INT-MULTI-SSBO

To get a $O(m)$ -approximation we follow the same approach as in Section 5.1. For each i we solve the restriction of the MULTI-SSBO problem on the i th scenario, *i.e.*, the quadratic program **(Q4)** as shown in Fig. 5, and then take the best of these solutions. It is easy to see that an optimal solution of **(Q4)** satisfies $\alpha_i = \min \left\{ 1, \frac{B_i}{\sum_{j=1}^n \sum_{k=1}^s w_{i,j,k}} \right\}$. For any fixed value of α_i , **(Q4)** is known in the literature as the *multiple-choice* Knapsack problem with sn objects divided into n classes and a knapsack capacity of B_i/α_i ; a $O(1)$ -approximation algorithm for this problem that runs in $O(ns^2)$ time is known [19].

We next show that algorithms for the single-slot case can be used for the multi-slot model with appropriate multiplicative factors in the approximation ratio.

Lemma 3. *There exists a $O(s\kappa \log \Delta)$ -approximation (respectively, $O(s \log^2(m+n)\kappa \log \Delta)$ -approximation) algorithm for FRAC-MULTI-SSBO (respectively, INT-MULTI-SSBO).*

Proof. We first prove our claim for For FRAC-MULTI-SSBO. Consider the quadratic program **(Q2)'** obtained from the quadratic program **(Q2)** for FRAC-MULTI-SSBO by removing the constraints $\sum_{k=1}^s x_{j,k} \leq 1$ for $1 \leq j \leq n$. If OPT and OPT' are the optimal values of the objective functions of **(Q2)** and **(Q2)'**, respectively, then obviously $\text{OPT}' \geq \text{OPT}$. A straightforward inspection shows that **(Q2)'** can be written down in the same form as **(Q1)** with sn variables and m constraints. Thus, using the already proven result of Theorem 1(i) we obtain a solution for **(Q2)'** whose objective value is $\frac{\text{OPT}'}{\kappa \log(\max_{j,k} d_{j,k})} = \frac{\text{OPT}'}{\kappa \log \Delta} \geq \frac{\text{OPT}}{\kappa \log \Delta}$. To convert this to a solution of FRAC-MULTI-SSBO (*i.e.*, to satisfy the constraints $\sum_{k=1}^s x_{j,k} \leq 1$ for each j) we divide each $x_{j,k}$ by $\sum_{k=1}^s x_{j,k}$ which decreases the total payoff by no more than a factor of s .

The result for INT-MULTI-SSBO follows by translating the above worst-case approximation bound for FRAC-MULTI-SSBO to a worst-case approximation of INT-MULTI-SSBO via the following lemma.

Lemma 4. (APPROXIMATING INT-MULTI-SSBO VIA FRAC-MULTI-SSBO) *Suppose that we have a η -approximation for FRAC-MULTI-SSBO. Then, we also have a $O(\eta\gamma)$ approximation for INT-MULTI-SSBO where $\gamma = \begin{cases} \log m, & \text{if } s = 1 \\ \log^2(m+n), & \text{otherwise} \end{cases}$*

Proof. For a particular value of the vector $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_m)$, **(Q2)** reduces to a linear program on the variables $\mathbf{x} = (x_1, x_2, \dots, x_n)$. For ease of description, we consider the case of $s = 1$ first (*i.e.*, the case of FRAC-SSBO). An inspection of **(Q1)** reveals that this linear program has exactly n variables and m inequalities, where the i th inequality D_i (for $1 \leq i \leq m$) is of the form:

$$D_i \stackrel{\text{def}}{=} \alpha_i \left(\sum_{j=1}^n w_{i,j} x_j \right) \leq B_i$$

Consider a solution $\mathbf{x}^f = (x_1^f, x_2^f, \dots, x_n^f)$. and $\alpha^f = (\alpha_1^f, \alpha_2^f, \dots, \alpha_m^f)$, of FRAC-SSBO with $\mathcal{L} = \sum_{i=1}^m \sum_{j=1}^n \alpha_i^f y_{i,j} x_j^f$ as the value of its objective. We may assume that $\mathcal{L} > 100 \ln m$ since otherwise the approximation guarantee can be trivially achieved. We employ the following randomized rounding scheme to transform this solution to a solution of INT-SSBO:

(* Quadratic program (Q4) *)

maximize $\alpha_i \left(\sum_{j=1}^n \sum_{k=1}^s y_{i,j,k} x_{j,k} \right)$

subject to

$$\alpha_i \left(\sum_{j=1}^n \sum_{k=1}^s w_{i,j,k} x_{j,k} \right) \leq B_i$$

$$\forall 1 \leq j \leq n: \sum_{k=1}^s x_{j,k} \leq 1$$

$$0 \leq \alpha_i \leq 1$$

$$\forall 1 \leq j \leq n \forall 1 \leq k \leq s: 0 \leq x_{j,k} \leq 1$$

Figure 5: MULTI-SSBO restricted to the i th scenario.

- For $i = 1, 2, \dots, n$, we round x_i^f randomly to 0 and 1 with probabilities x_i^f and $1 - x_i^f$, respectively. Let $x_i \in \{0, 1\}$ be the resulting random variable.
- We return $\mathbf{x} = (x_1, x_2, \dots, x_n)$ and $\boldsymbol{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_m)$ as our solution where $\alpha_i = \frac{\alpha_i^f}{100 \ln m}$ for $1 \leq i \leq m$.

Let $\mathcal{L}' = \sum_{i=1}^m \sum_{j=1}^n \alpha_i y_{i,j} x_j$ be the new value of the objective and let \mathcal{E}_i be the event that inequality D_i holds for this randomized solution. By linearity of expectation $\mathbb{E}[\mathcal{L}'] = \frac{\mathcal{L}}{100 \ln m}$. Consider the inequality D_i , and let $\alpha'_i = \frac{\alpha_i}{B_i + 1}$. By linearity of expectation,

$$\mathbb{E} \left[\alpha'_i \left(\sum_{j=1}^n w_{i,j} x_j \right) \right] = \frac{1}{100 \ln m} \times \frac{1}{B_i + 1} \times \alpha_i^f \left(\sum_{j=1}^n w_{i,j} x_j^f \right) < \frac{1}{100 \ln m} \times \frac{B_i}{B_i + 1}$$

Each random value $\alpha'_i w_{i,j} x_j$ can be thought of as an independent Poisson trial with a probability of success (*i.e.*, a value of 1) as $\alpha'_i w_{i,j} x_j$. Thus, using standard Chernoff bound [20, Exercise 4.1], we get:

$$\Pr[\mathcal{E}_i \text{ does not hold}] = \Pr \left[\alpha_i \left(\sum_{j=1}^n w_{i,j} x_j \right) > B_i \right] = \Pr \left[\alpha'_i \left(\sum_{j=1}^n w_{i,j} x_j \right) > \frac{B_i}{B_i + 1} \right] < e^{-3 \ln m} < \frac{1}{m^2}$$

In a similar manner, one can show that $\Pr[\mathcal{L}' < \frac{\mathcal{L}}{200 \ln m}] < \frac{1}{m}$. Thus, finally, using union bounds, we get

$$\Pr \left[\mathcal{L}' \geq \frac{\mathcal{L}}{200 \ln m} \bigwedge (\bigwedge_{i=1}^m \mathcal{E}_i \text{ holds}) \right] \geq 1 - \Pr \left[\mathcal{L}' < \frac{\mathcal{L}}{200 \ln m} \right] - \left(\sum_{i=1}^m \Pr[\mathcal{E}_i \text{ does not hold}] \right) > 1 - \frac{2}{m}$$

Thus, we achieve the desired approximation bounds with $1 - o(1)$ probability.

For the case of $s > 1$ (*i.e.*, FRAC-MULTI-SSBO), the same approach with some modifications work. In a nutshell, we have n additional constraints F_j (for $j = 1, 2, \dots, n$) of the form $\sum_{k=1}^s x_{j,k} \leq 1$. Thus, the total number of inequalities/equalities is $m + n$ and we need to do the analysis with “ $\ln(n + m)$ ” replacing “ $\ln m$ ”. The only additional part that needs to be done is to show how to handle the F_j constraints. Notice that the set of variables involved in F_j are disjoint from the set of variables in any other $F_{j'}$ for $j' \neq j$. After rounding, we have $\sum_{k=1}^s x_{j,k} \leq 100 \ln(m + n)$. We now select one of these variables x_{j_1} to $x_{j,s}$, say $x_{j,\ell}$, such that $x_{j,\ell} = \max_{1 \leq k \leq s} \{ \sum_{i=1}^m \alpha_i x_{j,k} y_{i,j,k} \}$, set $x_{j,\ell} = 1$ and set $x_{j,k} = 0$ for $k \neq \ell$. After all these normalizations, we loose an additional factor of $100 \ln(m + n)$ and all constraints are satisfied. \square

Note that the claim in Lemma 4 is “pessimistic” in nature; indeed, as our claim in Theorem 1 shows, for arbitrary parameter range both INT-SSBO and FRAC-SSBO can be approximated to within the same ratio. \square

6 Approximation-hardness Results for SSBO and MULTI-SSBO (main result (R2))

6.1 Approximation-hardness Bounds for SSBO

Theorem 5 (Logarithmic inapproximability). *There exists instances of INT-SSBO and FRAC-SSBO, with n keywords and $m = n$ scenarios each with equal probability, such that, unless $\text{ZPP} = \text{NP}$, any*

polynomial-time algorithm for solving these problems must have an approximation ratio of any one of the following:

- $\Omega(m^{1-\varepsilon})$ (and, thus, also $\Omega(n^{1-\varepsilon})$), or
- $\Omega(\kappa \log^{1-\varepsilon} d_n)$.

where $0 < \varepsilon < 1$ is any constant.

Proof. We construct instances of SSBO with n keywords and $m = n$ scenarios such that, for¹⁶ any κ and any values of $c_{i,j}$ in the range $[d_j, \kappa d_j]$, the claimed lower bound holds. We use the reformulation of FRAC-SSBO and INT-SSBO as a bipartite quadratic program **(Q2)** as discussed in Section 4.

The standard maximum independent set (MIS) problem is defined as follows. We are given an undirected graph $G = (V, E)$. A subset of vertices $V' \subseteq V$ is called *independent* if for *every* two vertices $u, v \in V'$ we have $\{u, v\} \notin E$. The goal is to find an independent subset of vertices of *maximum* cardinality. It is known that MIS cannot be approximated to within a factor of $|V|^{1-\varepsilon}$ for any constant $0 < \varepsilon < 1$ unless $\text{ZPP} = \text{NP}$ [15].

For notational simplicity, let $n = |V|$ and $a = n^{12}$. Set $m = n$. Select an arbitrary order v_1, v_2, \dots, v_n of the vertices in V . Intuitively, the i th column and the $(n+1-i)$ th row of Y correspond to the vertex v_i and the entries of the matrix Y are such that they are 0 above the reverse diagonal and encodes the adjacency of vertices of G on or below the reverse diagonal. Formally,

$$y_{i,j} = \begin{cases} 0 & \text{if } i+j < n+1 \\ 1 & \text{if } i+j = n+1 \\ 1 & \text{if } i+j > n+1 \text{ and } \{v_{n-i+1}, v_j\} \in E \\ 0 & \text{if } i+j > n+1 \text{ and } \{v_{n-i+1}, v_j\} \notin E \end{cases}$$

Fix d_1, d_2, \dots, d_n as $d_1 = 1$ and $d_i = a d_{i-1}$ for $1 < i \leq n$. Thus, $\frac{c_{i,j_1}}{c_{i,j_2}} \geq \frac{d_{j_1}}{\kappa d_{j_2}} > n^6$ if $j_1 > j_2$. Let $B_i = c_{i,n+1-i}$ for $1 \leq i \leq m = n$. Remembering that $w_{i,j} = c_{i,j} y_{i,j}$ for all i and j , we have:

$$w_{i,j} = \begin{cases} 0 & \text{if } i+j < n+1 \\ & \text{or if } i+j > n+1 \text{ and } \{v_{n-i+1}, v_j\} \notin E \\ c_{i,j} & \text{if } i+j = n+1 \\ & \text{or if } i+j > n+1 \text{ and } \{v_{n-i+1}, v_j\} \in E \end{cases}$$

Note that $n^{1-\varepsilon} = m^{1-\varepsilon} = \Omega(\kappa \log^{1-\varepsilon} d_n)$ since $d_n = n^{12n}$ and $\kappa = \text{poly}(\log(m+n)) = \text{poly}(\log(n))$. Let Δ_{ind} and Δ_{Q1} be the maximum number of independent vertices in G and an optimal value of the objective of the fractional or integral version of **(Q1)**, respectively.

Lemma 6. $\Delta_{\text{Q1}} \geq \Delta_{\text{ind}}$.

Proof. Consider an optimal solution V' of MIS on G with $|V'| = \Delta_{\text{ind}}$. We generate a solution of **(Q2)** by setting

$$x_i = \alpha_{n-i+1} = \begin{cases} 1, & \text{if } v_i \in V' \\ 0, & \text{otherwise} \end{cases}$$

Note that, since V' is an independent set, if $i+j > n+1$, $v_i \in V'$ and $\{v_i, v_j\} \in E$ then $v_j \notin V'$ and thus $x_i = \alpha_{n-i+1} = 1$ and $x_j = \alpha_{n-j+1} = 0$.

¹⁶Remember that in Section 2.5 we fixed bounds on κ , namely, $\kappa = O(\text{poly}(\log(m+n)))$.

First, we show that this is indeed a valid solution of **(Q1)**. For any $1 \leq i \leq n-1$, consider the constraint

$$\alpha_{n-i+1} \left(\sum_{j=1}^n w_{n-i+1,j} x_j \right) \leq B_{n-i+1}.$$

If $\alpha_{n-i+1} = 0$, then the constraint is obviously satisfied since $B_{n-i+1} > 0$. Otherwise, $\alpha_{n-i+1} = x_i = 1$ and thus,

$$\alpha_{n-i+1} \left(\sum_{j=1}^n w_{n-i+1,j} x_j \right) = \sum_{j=1}^n w_{n-i+1,j} x_j = c_i + \sum_{\substack{i+j > n+1 \\ \{v_i, v_j\} \in E}} w_{n-i+1,j} x_j = c_{n+1-i,i} = B_{n+1-i}$$

Thus, all the constraints are satisfied. Finally, the value of the objective function is

$$\sum_{i=1}^m \sum_{j=1}^n \alpha_i x_j y_{i,j} = \sum_{\substack{i+j=n+1 \\ v_j \in V'}} \alpha_i x_j = \sum_{v_j \in V'} x_j = \Delta_{\text{ind}}$$

and thus $\Delta_{Q1} \geq \Delta_{\text{ind}}$. □

For the other direction, we first need a normalization lemma.

Lemma 7 (Normalization lemma). *Consider an optimal solution of **(Q1)** with an objective value of Δ_{Q1} . Then, we can transform this solution to another solution of **(Q1)** of objective value Δ'_{Q1} such that:*

- (a) $x_i \in \{0, 1\}$ for each i ;
- (b) $\Delta'_{Q1} \geq \Delta_{Q1} - 1$; and
- (c) if $\{x_i, x_j\} \in E$ then $x_i + x_j \leq 1$.

Proof. Suppose that we are given an optimal solution of **(Q1)** with an objective value of Δ_{Q1} . First, we note some properties of this solution.

Proposition 3. *The following statements are true:*

- (i) for every i , $\alpha_{n-i+1} x_i \leq 1$, and
- (ii) for every i and j , if $i+j > n+1$ and $\{v_i, v_j\} \in E$ then $\alpha_j x_i \leq n^{-6}$.

Proof. Consider the constraint $\alpha_{n-i+1} \left(\sum_{j=1}^n w_{n-i+1,j} x_j \right) \leq B_{n-i+1} = c_{n-i+1,i}$.

Since $w_{n-i+1,i} = c_{n-i+1,i}$, (i) follows.

(ii) is equivalent to the claim that $\alpha_{n-i+1} x_j \leq n^{-6}$ if $j > i$. Since $\frac{c_{p,j}}{c_{q,i}} > n^6$ if $j > i$ (for any p and q), (ii) follows. □

Now we show how to “normalize” this solution such that each variable x_i is 0 or 1, and the total objective value does not decrease too much. Let $\Gamma = \sum_{i+j \neq n+1} \alpha_i x_j y_{i,j}$. By Proposition 3(ii), $\Gamma \leq n^2 \times n^{-6} = n^{-4}$. Thus, setting $\Phi = \sum_{i+j=n+1} \alpha_i x_j y_{i,j}$, it follows that $\Phi \leq \Delta_{Q1} \leq \Phi + n^{-4}$. Thus, subsequently we concentrate on the quantity Φ .

If $\alpha_{n-i+1} = 0$ for some i , then we can set $x_i = 0$ without changing the value of Φ . Let $I = \{n-i+1 \mid \alpha_{n-i+1} > 0 \text{ and } x_i > 0\}$. Consider the largest index $n-i+1 \in I$. There are two cases to consider:

Case 1: $x_i > n^{-3}$. By Proposition 3(i), $\alpha_{n-i+1} < n^{-3}$ and $\alpha_{n-j+1}x_j \leq \alpha_{n-j+1} < n^{-3}$ for every $j > i$ such that $\{v_i, v_j\} \in E$.

We set $\alpha_{n-i+1} = x_i = 1$ and set $x_j = \alpha_{n-j+1} = 0$ for every $j > i$ such that $\{v_i, v_j\} \in E$. The change in Φ is at most $n \times n^{-3} = n^{-2}$.

Case 2: $x_i \leq n^{-3}$. We set $\alpha_{n-i+1} = x_i = 0$. The change in Φ is at most n^{-3} .

We now remove the index $n - i + 1$ from I and continue with the next largest index. We continue until $I = \emptyset$. Since $|I| \leq n$, the total change in Φ is at most $n^{-1} < 1 - n^{-4}$.

To complete the proof, we select vertices v_j in the independent set if $x_j = 1$. \square

To finish the proof of Theorem 5, we simply select those vertices v_i for the independent set such that $x_i = 1$. We have now shown that $\Delta_{\text{ind}} \leq \Delta_{Q1} \leq \Delta_{\text{ind}} - 1$. Thus, since Δ_{ind} and Δ_{Q1} are within a constant factor of each other and Δ_{ind} cannot be approximated to within a factor of $n^{1-\varepsilon}$ for any constant $0 < \varepsilon < 1$, Δ_{Q1} cannot be approximated to within a factor of $\Omega(n^{1-\varepsilon})$, or $\Omega(m^{1-\varepsilon})$, or $\Omega(\kappa \log^{1-\varepsilon} d_n)$. \square

6.2 Approximation Hardness Results for MULTI-SSBO

A first natural approach for this would be to generalize the approximation hardness result for the single-slot case (Q1) in Theorem 5 to the multi-slot case (Q2). This can be trivially done by copying the construction of the single-slot case to *one* of the slots in the multi-slot case. However, after this, one can observe that:

the construction for the single-slot case cannot again be copied to another slot because of the constraints in Equation (2) which states that at most one selection variable in each slot can be set to 1.

Formally, the lower bound construction for (Q1) can be extended to (Q2) as follows:

- Identify $y_{i,j,1}$ of (Q2) with $y_{i,j}$ of (Q1) and set $y_{i,j,2} = y_{i,j,3} = \dots = y_{i,j,s} = 0$ in (Q2).
- Identify $c_{i,j,1}$ of (Q2) with $c_{i,j}$ of (Q1) and set $c_{i,j,2} = c_{i,j,3} = \dots = c_{i,j,s} = 0$ in (Q2).
- Identify $x_{j,k,1}$ of (Q2) with x_j of (Q1).

This leads to the following approximation hardness result.

Corollary 8. *There exists instances of INT-MULTI-SSBO and FRAC-MULTI-SSBO, with n keywords, $m = n$ scenarios each with equal probability and s slots, such that, unless $ZPP = NP$, any polynomial-time algorithm for solving these problems must have an approximation ratio of $\Omega(n^{1-\varepsilon})$ or $\Omega(\kappa \log^{1-\varepsilon} d_n)$, where $0 < \varepsilon < 1$ is any constant.*

However, note that the lower bound in the above corollary **does not involve** s . This is unfortunate, since one would expect that the asymptotic complexity of MULTI-SSBO problems should depend on the number of slots s , especially when s is large. So, a natural question to ask is:

can we obtain an approximation hardness result for MULTI-SSBO problems that shows dependencies on s ?

Our next result in this section shows that one can prove a “super poly-logarithmic but sub-polynomial” bound¹⁷ of $2^{\log^{1-\varepsilon}(ns)}$ that depends on s . The bound is not completely satisfactory because of its sub-exponential nature, but nonetheless provides evidence that the approximation hardness of MULTI-SSBO problems depends strongly on s . How do the two bounds $n^{1-\varepsilon}$ and $2^{\log^{1-\varepsilon}(ns)}$ compare? Note that:

- Since $2^{\log^{1-\varepsilon}(ns)} = o((ns)^\delta)$ for any constant $\delta > 0$, $n^{1-\varepsilon}$ dominates $2^{\log^{1-\varepsilon}(ns)}$ if s is not too large compared to n , *e.g.*, $s = O(n)$.
- On the other hand, $2^{\log^{1-\varepsilon}(ns)}$ dominates $n^{1-\varepsilon}$ if s is large compared to n , *e.g.*, $s = \Omega(n^{\log n})$.

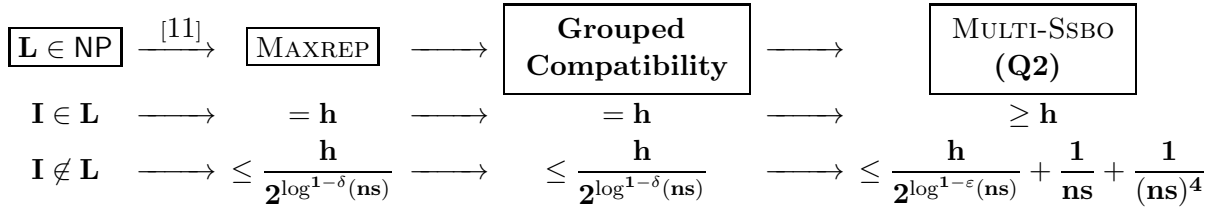
Thus, neither bound is subsumed by the other for all values of parameters.

Theorem 9 (Inapproximability of multi-slot case).

- (a) (DEPENDENCE ON NUMBER OF SLOTS) *There exists instances of INT-MULTI-SSBO and FRAC-MULTI-SSBO with n keywords, s slots and $\kappa = 1$ such that, unless $\text{NP} \subseteq \text{DTIME}(n^{\text{poly}(\log(n))})$, any polynomial-time algorithm for solving these problems must have an approximation ratio of $2^{\log^{1-\varepsilon}(ns)}$ where $0 < \varepsilon < 1$ is any constant.*
- (b) INT-MULTI-SSBO is MAX-SNP-hard for $s = 2$ even when $\kappa = 1$ and $c_{j,k} = 1$ for all j and k .

6.2.1 Proof of Theorem 9 (a)

A schematic diagram of the entire reduction is shown below.



Our beginning is a one-round two-prover system on which the parallel repetition approach applies. It would be more convenient to describe the problem in a graph-theoretic form as the MAXREP problem [18]. We are given a bipartite graph $\mathcal{G} = (\mathcal{A}, \mathcal{B}, E)$ with $|\mathcal{A}| = |\mathcal{B}|$. Also is given a partition of \mathcal{A} into n equal-size subsets $\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_n$, each with s elements, and a partition of \mathcal{B} into n equal-size subsets $\mathcal{B}_1, \mathcal{B}_2, \dots, \mathcal{B}_n$, each having s elements (and, thus, $|\mathcal{A}| = |\mathcal{B}| = ns$). These partitions define a natural “bipartite super-graph” \mathcal{H} in the following manner. \mathcal{H} has a “super-vertex” for every \mathcal{A}_i (the left partition) and a “super-vertex” for every \mathcal{B}_j (the right partition). There exists a “super-edge” between the super-vertex \mathcal{A}_i and the super-vertex \mathcal{B}_j if and only if there exists $u \in \mathcal{A}_i$ and $v \in \mathcal{B}_j$ such that $\{u, v\}$ is an edge of \mathcal{G} . It is also given that \mathcal{H} is a d -regular graph for some d . Thus, the number of super-edge h is given by $h = dn$. A pair of nodes u and v “witnesses” a super-edge $\{\mathcal{A}_i, \mathcal{B}_j\}$ provided $u \in \mathcal{A}_i$, $v \in \mathcal{B}_j$ and the edge $\{u, v\}$ exists in \mathcal{G} . A set of nodes S of \mathcal{G} witnesses a super-edge if and only if there exists at least one pair of nodes in S that witnesses the super-edge. In MAXREP, we are supposed to select a single vertex from each \mathcal{A}_i and a single vertex from each \mathcal{B}_j . The goal of MAXREP is to *maximize* the *size* of its solution, namely the number of super-edges witnessed. The following result follows from [11] (see also [18] for a self-contained description).

¹⁷ $2^{\log^{1-\varepsilon}a}$ is $\omega(\text{poly}(\log a))$ and $o(a^\delta)$ for any constant $\delta > 0$.

Theorem 10. [11, 18] *Let $L \in \text{NP}$ and $0 < \delta < 1$ be any fixed constant. Then, there exists a reduction running in quasi-polynomial time, namely in time $n^{\text{poly}(\log n)}$, that given an instance I of L produces an instance I' of MAXREP such that:*

- if $I \in L$ then I' has a solution of size h ;
- if $I \notin L$ then I' has every solution of size at most $\frac{h}{2^{\log^{1-\delta}(ns)}}$.

Thus, the above theorem provides a $2^{\log^{1-\delta}(ns)}$ -inapproximability for MAXREP under the complexity-theoretic assumption of $\text{NP} \not\subseteq \text{DTIME}(n^{\text{poly}(\log n)})$.

Let L be any language in NP. We first use the above theorem to translate an instance I of L to an instance I' of MAXREP as described. We next translate this instance I' to an instance of an intermediate problem, which we call as the Grouped Compatibility (GROUP-COMPATIBILITY) problem, that we define next.

For notational convenience, we define two mappings σ and π such that $\sigma(u) = \mathcal{A}_i$ and $\pi(v) = \mathcal{B}_j$ if $u \in \mathcal{A}_i$ and $v \in \mathcal{B}_j$, respectively. The GROUP-COMPATIBILITY problem is derived from MAXREP as follows. Given an instance of MAXREP as described above, we construct a $\mathcal{G}' = (\mathcal{V}', E')$ in the following manner. \mathcal{V}' has a group $\mathcal{V}'_{i,j}$ of s^2 vertices for every pair \mathcal{A}_i and \mathcal{B}_j ; thus there are n^2 such groups. There is a vertex u_v in $\mathcal{V}'_{i,j}$ for every pair of vertices $u \in \mathcal{A}_i$ and $v \in \mathcal{B}_j$. The weight $w(u_v)$ of such a vertex u_v is set as follows: if $\{u, v\} \in E$ then $w(u_v) = 1$ else $w(u_v) = 0$. The edges of \mathcal{G}' are defined as follows. For two vertices u_v and $u'_{v'}$, the edge $\{u_v, u'_{v'}\}$ exists in \mathcal{G}' if and only if *exactly* one of the following two conditions (★) and (★★) are satisfied:

(★) $\sigma(u) = \sigma(u')$ and $u \neq u'$;

(★★) $\pi(v) = \pi(v')$ and $v \neq v'$.

The goal of GROUP-COMPATIBILITY is to select *at most one* vertex from every group $\mathcal{V}'_{i,j}$ such that no two selected vertices are adjacent and the *size* of the solution, namely the sum of weights of selected vertices, is *maximized*. One could translate a lower bound for MAXREP to a lower bound for GROUP-COMPATIBILITY preserving inapproximability in the following manner.

Lemma 11. *MAXREP has an optimal solution of size y if and only if GROUP-COMPATIBILITY has an optimal solution of size y .*

Proof. Consider a solution of MAXREP of size y . We select the set of vertices u_v as a solution of GS for every pair of vertices $u \in \mathcal{A}$ and $v \in \mathcal{B}$ selected by MAXREP. That this is indeed a valid solution for GROUP-COMPATIBILITY can be seen as follows. Since MAXREP selects exactly one vertex from \mathcal{A}_i and every \mathcal{B}_j , for any two selected vertices u_v and $u'_{v'}$ if $\sigma(u) = \sigma(u')$ then $u = u'$ or if $\pi(v) = \pi(v')$ then $v = v'$, and moreover if MAXREP selects $u \in \mathcal{A}_i$ and $v \in \mathcal{B}_j$ then only the vertex u_v is selected from the group $\mathcal{V}'_{i,j}$. For each pair $u \in \mathcal{A}_i$ and $v \in \mathcal{B}_j$ selected by MAXREP, the pair contributes 1 to the size of MAXREP if and only if $\{u, v\} \in E$ which holds if and only if $w(u_v) = 1$. Thus, the size of the MAXREP solution is identical to that of GROUP-COMPATIBILITY.

For the other direction, consider a solution of size y of GROUP-COMPATIBILITY. For every vertex u_v in the solution of GROUP-COMPATIBILITY, we select the vertices $u \in \mathcal{A}$ and $v \in \mathcal{B}$ in the solution of MAXREP. Now, we note the following. If GROUP-COMPATIBILITY selects two vertices u_v and $u'_{v'}$ then either

- $\sigma(u) \neq \sigma(u')$ and $\pi(v) \neq \pi(v')$; or
- $\sigma(u) = \sigma(u')$ and $u = u'$; or

- $\pi(v) = \pi(v')$ and $v = v'$.

In all cases, at most one vertex from each partition of \mathcal{A} or \mathcal{B} are selected.

Now, suppose that the above solution of MAXREP did not select a vertex in some partition of $\mathcal{A} \cup \mathcal{B}$, say \mathcal{A}_i (the case of \mathcal{B}_j is similar). This means the solution of GROUP-COMPATIBILITY did not include a vertex u_v for any $u \in \mathcal{A}_i$ and any $v \in \mathcal{B}$. Pick any $v \in \mathcal{B}$ such that x_v is in the solution of GROUP-COMPATIBILITY for some $x \in \mathcal{A}$ (if GROUP-COMPATIBILITY does not contain any such v pick a $v \in \mathcal{B}$ arbitrarily). Pick any $u \in \mathcal{A}_i$ and suppose that we add this vertex u_v to the solution of GROUP-COMPATIBILITY. It is easy to see that u_v is not adjacent to any other vertex in the solution of GROUP-COMPATIBILITY thereby increasing the size of the solution of GROUP-COMPATIBILITY by 1 and contradicting the optimality of the solution of GROUP-COMPATIBILITY.

Finally, if $w(u_v) = 1$ for a vertex u_v in the solution of GROUP-COMPATIBILITY then both vertices u and v are added to the solution of MAXREP and since $\{u, v\} \in E$ the pair of vertices (u, v) witnesses one more super-edge. Thus, the size of the solution of MAXREP is precisely y . \square

For any constant $0 < \delta < 1$, let $0 < \varepsilon < 1$ be the corresponding constant such that $(\log_2(n^2 s^2))^{1-\varepsilon} = (\log_2(ns))^{1-\delta}$. Combining Lemma 11 and Theorem 10, the following lemma follows.

Lemma 12. *Let $L \in \text{NP}$ and $0 < \varepsilon < 1$ be any fixed constant. Then, there exists a reduction running in quasi-polynomial time, namely in time $n^{\text{poly}(\log n)}$, that given an instance I of L produces an instance I' of GROUP-COMPATIBILITY with n groups of vertices, each group having s vertices, such that, for some h :*

- if $I \in L$ then I' has a solution of size h ;
- if $I \notin L$ then I' has every solution of size at most $h/2^{\log^{1-\varepsilon}(ns)}$.

Now, we provide a translation of the instance I' of GROUP-COMPATIBILITY as to an instance I'' of MULTI-SSBO using the quadratic programming reformulation (Q2). Let $\mathcal{G} = (\mathcal{V}, E)$ be the instance I' of GROUP-COMPATIBILITY with a given partition of \mathcal{V} into n groups of vertices $\mathcal{V}_1, \mathcal{V}_2, \dots, \mathcal{V}_n$ where $|\mathcal{V}_i| = s$. Let $u_{i,1}, u_{i,2}, \dots, u_{i,s}$ be an arbitrary ordering of the s vertices of \mathcal{V}_i and let v_1, v_2, \dots, v_{ns} denote the ordering

$$u_{1,1}, u_{1,2}, \dots, u_{1,s}, u_{2,1}, u_{2,2}, \dots, u_{2,s}, \dots, u_{n,1}, u_{n,2}, \dots, u_{n,s}$$

of vertices. We construct an instance I'' of MULTI-SSBO in the following manner. For every \mathcal{V}_i , there is a keyword K_i with s slots with $c_{i,j} = (ns)^{6(i-1)s+6j}$. Intuitively, we associate the j th slot of K_i with the vertex $u_{i,j}$ and scenarios 1 through $m = ns$ correspond to an enumeration of vertices of \mathcal{G} in the order $v_{ns}, v_{ns-1}, \dots, v_1$. The entries $y_{i,j,k}$ of the profit matrix Y and the quantities B_1, B_2, \dots, B_{ns} in (Q2) are specified as follows. Let $ns - i + 1 = (j' - 1)s + k'$ for some integers $1 \leq j' \leq n$ and $1 \leq k' \leq s$. Then,

$$y_{i,j,k} = \begin{cases} 0 & \text{if } j < j' \\ 0 & \text{if } j = j' \text{ and } k < k' \\ w(u_{j,k}) & \text{if } j = j' \text{ and } k = k' \\ 1 & \text{if } j = j' \text{ and } k > k' \\ 1 & \text{if } j > j' \text{ and } \{u_{j',k'}, u_{j,k}\} \in E \\ 0 & \text{if } j > j' \text{ and } \{u_{j',k'}, u_{j,k}\} \notin E \end{cases}$$

and $B_i = c_{j',k'} = (ns)^{6(j'-1)s+6k'}$. It can be shown that the reduction produces the desired inapproximability by proving the following lemma and noting that $(ns)^{-1} + (ns)^{-4} < 1$.

Lemma 13. *The following two statements are true.*

- (a) *If GROUP-COMPATIBILITY has a solution of size h , then MULTI-SSBO has a solution of size at least h .*
- (b) *If GROUP-COMPATIBILITY has no solution of size more than $h/2^{\log^{1-\varepsilon}(ns)}$, then MULTI-SSBO has no solution of size more than $\frac{h}{2^{\log^{1-\varepsilon}(ns)}} + \frac{1}{ns} + \frac{1}{(ns)^4}$.*

Proof.

(a) Consider a solution of I' of size h and *renumber* the indices, without loss of generality, such that if a vertex is picked from the partition \mathcal{V}_j it is vertex $u_{j,1}$. Let S be the set of vertices selected in the solution of I' . Now, for every j such that the vertex $u_{j,1} \in S$, we set $x_{j,1} = \alpha_{ns-(j-1)s} = 1$.

We first calculate the total payoff, i.e., the objective value of **(Q2)**, obtained by this solution. Consider the quantity $\text{payoff}_i = \alpha_i \left(\sum_{j=1}^n \sum_{k=1}^s x_{j,k} y_{i,j,k} \right)$. If $i \neq ns - (j' - 1)s$ for some $1 \leq j' \leq s$, then $\alpha_i = 0$ and thus $\text{payoff}_i = 0$. Now suppose that $i = ns - (j' - 1)s$; thus $ns - i + 1 = (j' - 1)s + 1$. Since S is a solution of GROUP-COMPATIBILITY, if $x_{j,1}, x_{j',1} \in S$ then $\{u_{j,1}, u_{j',1}\} \notin E$. Thus, $\text{payoff}_i = \alpha_i x_{j',1} y_{i,j',1} = w(u_{j',1})$. Consequently, $\text{payoff} = \sum_i \text{payoff}_i = \sum_{u_{j',1} \in S} w(u_{j',1}) = h$.

Next, we show that all the constraints of **(Q2)** are satisfied. The constraint $\sum_{k=1}^s x_{j,k} \leq 1$ is obviously satisfied for every j . So, we consider, for an arbitrary i , the constraint

$$\alpha_i \left(\sum_{j=1}^n \sum_{k=1}^s y_{i,j,k} c_{j,k} x_{j,k} \right) \leq B_i.$$

If $i \neq ns - (j' - 1)s$ for some $1 \leq j' \leq s$, then $\alpha_i = 0$ and the constraint is obviously satisfied. So, assume that $i = ns - (j' - 1)s$ for some $1 \leq j' \leq s$. Then,

$$\alpha_i \left(\sum_{j=1}^n \sum_{k=1}^s y_{i,j,k} c_{j,k} x_{j,k} \right) = w(u_{j',1}) c_{j',1} \leq c_{j',1} = B_i$$

(b) We need to show that a “small-size” solution of I' must imply a “similar small-size” solution of I'' . Since a solution of I'' allows fractional values of $x_{j,k}$ ’s, we first need a few normalization procedures which changes the objective value of any solution of **(Q2)** for the instance I'' only by a $o(1)$ additive factor.

Lemma 14. *Consider any $1 \leq i \leq ns$ and let $ns - i + 1 = (j' - 1)s + k'$ for some $1 \leq j' \leq n$ and $1 \leq k' \leq s$. Then, $\alpha_i y_{i,j,k} x_{j,k} \leq \begin{cases} w(u_{j,k}) & \text{if } j = j' \text{ and } k = k' \\ (ns)^{-6} & \text{otherwise} \end{cases}$.*

Proof. Suppose that $j = j'$ and $k = k'$; thus $y_{i,j,k} = w(u_{j,k})$. If $w(u_{j,k}) = 0$ then $\alpha_i y_{i,j,k} x_{j,k} = 0 = w(u_{j,k})$. If $w(u_{j,k}) = 1$ then, because of the constraint $\alpha_i \left(\sum_{j=1}^n \sum_{k=1}^s y_{i,j,k} c_{j,k} x_{j,k} \right) \leq B_i$, we must have $\alpha_i y_{i,j,k} c_{j,k} x_{j,k} \leq B_i$. Thus, $\alpha_i y_{i,j,k} x_{j,k} \leq \frac{B_i}{c_{j,k}} = 1 = w(u_{j,k})$.

If $j < j'$ or if $j = j'$ and $k < k'$ or if $j > j'$ and $\{u_{j,k}, u_{j',k'}\} \notin E$ then $y_{i,j,k} = 0 < (ns)^{-6}$.

If $j > j'$ and $\{u_{j,k}, u_{j',k'}\} \in E$ then, because of the constraint $\alpha_i \left(\sum_{j=1}^n \sum_{k=1}^s y_{i,j,k} c_{j,k} x_{j,k} \right) \leq B_i$, we must have $\alpha_i y_{i,j,k} c_{j,k} x_{j,k} \leq B_i$. Thus,

$$\alpha_i y_{i,j,k} x_{j,k} \leq \frac{B_i}{c_{j,k}} = \frac{c_{j',k'}}{c_{j,k}} = (ns)^{6(j'-1)s+6k'-6(j-1)s-6k} = (ns)^{6(j'-j)s+6(k'-k)} \leq \frac{1}{(ns)^6}$$

□

Let $\Delta = \sum_{\substack{i,j,k: \\ ns-i+1 \neq (j-1)s+k}} \alpha_i x_{j,k} y_{i,j,k}$ and $\Gamma = \sum_{\substack{i,j,k: \\ ns-i+1 = (j-1)s+k}} \alpha_i x_{j,k} y_{i,j,k}$; thus

$$\text{payoff} = \sum_{i=1}^m \text{payoff}_i = \alpha_i \left(\sum_{j=1}^n \sum_{k=1}^s y_{i,j,k} x_{j,k} \right) = \Gamma + \Delta.$$

We first show that $\Delta = o(1)$.

Lemma 15. $\Delta \leq (ns)^{-4}$.

Proof. Since $1 \leq i \leq ns$, $1 \leq j \leq n$ and $1 \leq k \leq s$, using Lemma 14 we get $\Delta \leq (ns)^2 / (ns)^6 = (ns)^{-4}$. \square

Thus, we now concentrate on the quantity Γ .

Lemma 16 (Normalization lemma). *Given a solution of I'' of MULTI-SSBO, it is possible to alter the solution in polynomial time such that if Γ' is the new value of Γ after the alteration then:*

(normalization within each slot) for each $1 \leq j \leq n$ $x_{j,k} \in \{0, 1\}$ and $\sum_{k=1}^s x_{j,k} \leq 1$;

(compatibility between slots) if $x_{j,k} = x_{j',k'} = 1$ then $\{\mathcal{V}_{j,k}, \mathcal{V}_{j',k'}\} \notin E$;

(small loss in payoff) $\Gamma' \geq \Gamma - (ns)^{-1}$.

Proof. In the sequel by triplets of indices i , j and k we mean a triplet such that $ns - i + 1 = (j - 1)s + k$. If $\alpha_i = 0$ or $y_{i,j,k} = w(u_{i,j}) = 0$ then we can set $x_{j,k} = 0$ without changing the value of Γ , so we assume that this is not the case.

Let j be the largest index such that $0 < x_{j,k} < 1$. Among all such k 's let k be the largest index such that $0 < x_{j,k} < 1$. If there is no such pairs of indices, our normalization procedure ends. Otherwise, by Lemma 14, $\alpha_i y_{i,j,k} x_{j,k} = \alpha_i x_{j,k} \leq 1$. We have the following cases.

- If $x_{j,k} \leq (ns)^{-3}$, we set $x_{j,k} = 0$. The resulting change in Γ is at most $(ns)^{-3}$.
- Otherwise $x_{j,k} > (ns)^{-3}$. By Lemma 14, we have:
 - (a) for $k' < k$ and i' such that $ns - i' + 1 = (j - 1)s + k'$, we have $\alpha_{i'} y_{i',j,k'} x_{j,k} = \alpha_{i'} x_{j,k} \leq (ns)^{-6}$. Since $x_{j,k} > (ns)^{-3}$ this implies $\alpha_{i'} \leq \alpha_{i'} x_{j,k} \leq (ns)^{-3}$. Thus, if we set $x_{j,k'} = 0$, the change in Γ due to this is at most $\alpha_{i'} y_{i',j,k'} x_{j,k'} \leq \alpha_{i'} \leq (ns)^{-3}$.
For each $k' < k$ we set $x_{j,k'} = 0$. The net change in Γ is at most $(ns)^{-2}$. Notice that after this change $x_{j,k} = 1$ and $x_{j,k'} = 0$ for every $k' \neq k$.
 - (b) For $j' < j$ and k' such that $\{u_{j,k}, u_{j',k'}\} \in E$, let i' be such that $ns - i' + 1 = (j' - 1)s + k'$. Then, we have $\alpha_{i'} y_{i',j',k'} x_{j,k} = \alpha_{i'} x_{j,k} \leq (ns)^{-6}$. Since $x_{j,k} > (ns)^{-3}$ this implies $\alpha_{i'} \leq \alpha_{i'} x_{j,k} \leq (ns)^{-3}$. Thus, if we set $x_{j',k'} = 0$, the change in Γ due to this is at most $\alpha_{i'} y_{i',j',k'} x_{j',k'} \leq \alpha_{i'} \leq (ns)^{-3}$.
For each such j' and k' as described above we set $x_{j',k'} = 0$. The net change in Γ is at most $(ns)^{-2}$. Notice that after this change if $x_{j',k'} > 0$ then $\{u_{j,k}, u_{j',k'}\} \notin E$. \square

We repeat the normalization procedure until it cannot be applied anymore. Note that both Step (a) and Step (b) above executes at most n times. Thus, the total change in Γ is at most $(ns)^{-1}$. To complete the proof, we select vertices $v_{j,k}$ in our solution if $x_{j,k} = 1$. \square

6.2.2 Proof of Theorem 9 (b)

We reduce the MAX-2SAT-5 problem¹⁸ to our problem. MAX-2SAT-5 is defined as follows. We are given a collection of m clauses C_1, C_2, \dots, C_m over n Boolean variables z_1, z_2, \dots, z_n , where every clause is a disjunction of exactly two literals and every variable occurs exactly 5 times (and, thus, $m = 5n/2$). The goal is to find an assignment of truth values to variables to satisfy a maximum number of clauses. This problem was shown to be MAX-SNP-hard in [10].

Given an instance of MAX-2SAT-5 we create an instance of INT-MULTI-SSBO (*i.e.*, **(Q2)**) with $s = 2$ as follows. Every variable z_j corresponds to a keyword \mathcal{K}_j with two slots. The variables $x_{j,1}$ and $x_{j,2}$ encode the truth assignments of the variable z_j with $x_{j,1} = 1$ indicating that z_j is true and $x_{j,2} = 1$ indicating that z_j is false; we will say that $x_{j,1}$ and $x_{j,2}$ are the slots corresponding to the literals z_j and $\neg z_j$, respectively. There are exactly m scenarios, each with probability $\frac{1}{m}$, defined in the following manner:

- $B_i = 1$ for $1 \leq i \leq m$.
- $c_{j,k} = 1$ for $1 \leq j \leq n$ and $1 \leq k \leq 2 = s$.
- For the i th clause C_i containing two *literals*, we have the i th scenario of the following form. Let $x_{j,k}$ and $x_{j',k'}$ be the slots corresponding to the two literals of the clause. Then we set $y_{i,j,k} = y_{i,j',k'} = 1$, and $y_{i,j,k} = 0$ if $j \neq j'$ or $k \neq k'$. For example, if $C_i = z_2 \vee (\neg z_3)$ then $y_{i,2,1} = y_{i,3,2} = 1$ and $y_{i,j,k} = 0$ for all other j and k .

An inspection of the construction reveals that it satisfies the following:

- Because this is an instance of INT-MULTI-SSBO, by Equation (2), for every $1 \leq j \leq n$, either $x_{j,1} = 1$ or $x_{j,2} = 1$ but *not both*. On the other hand, it is always possible to set at least one of the two variables $x_{j,1} = 1$ or $x_{j,2} = 1$ without increasing the total payoff. Thus setting these variables correspond to a truth assignment.
- A scenario contributes a payoff of 1 if and only if *at least* one of two slots have been selected. Thus, contribution of a scenario correspond to satisfying a clause.

By the above observations, we satisfy m' clauses if and only if the above instance of INT-MULTI-SSBO has a total payoff of m' .

7 Other Results

7.1 Improved Algorithms for Special Cases of SSBO and MULTI-SSBO

By the phrase “within an additive error of δ ” in Lemma 17 we mean that if our solution returns an objective value of x when the optimal value is y then $|x - y| \leq \delta$.

Lemma 17.

- (a) **(Fixed number of scenarios)** *If m is fixed, FRAC-MULTI-SSBO admits a pseudo-polynomial time solution with an absolute error of δ for any fixed $\delta > 0$, INT-SSBO admits a pseudo-polynomial time $O(1)$ -approximation and INT-MULTI-SSBO admits a pseudo-polynomial time $O(\log^2 n)$ -approximation.*

¹⁸Our reduction approach should also work if we start with MAX-2SAT- k for any constant k .

- (b) **(Fixed number of keywords)** If ns is fixed, then FRAC-MULTI-SSBO admits a **pseudo-polynomial** time solution with an absolute error of δ for any fixed $\delta > 0$.
- (c) **(Logarithmic number of keywords)** if $ns = O(\log m)$ then INT-MULTI-SSBO admits a polynomial time exact solution.
- (d) **(Fixed number of scenarios and polynomial bids)** If m is fixed and the maximum size of all the numbers, namely $\max \left\{ \max_{i,j,k} \{y_{i,j,k}\}, \max_i \{B_i\}, \max_i \left\{ \frac{1}{\varepsilon_i} \right\}, \max_{i,j,k} \{c_{i,j,k}\} \right\}$, is at most $\text{poly}(n)$ then INT-MULTI-SSBO admits a polynomial time solution with an absolute error of δ for any fixed $\delta > 0$.

Proof.

(a) and (b) We prove part (a) as follows (the proof for part (b) is similar). Consider the FRAC-MULTI-SSBO problem; let $y = \max_{\substack{1 \leq i \leq m \\ 1 \leq j \leq n \\ 1 \leq k \leq s}} \{y_{i,j,k}\}$.

Proposition 4. Let $\alpha^* = (\alpha_1^*, \alpha_2^*, \dots, \alpha_m^*)$ and $\mathbf{x}^* = (x_{1,1}^*, \dots, x_{1,s}^*, x_{2,1}^*, \dots, x_{2,s}^*, \dots, x_{n,1}^*, \dots, x_{n,s}^*)$ be the solution vectors for an optimal solution of value $\mathbb{E}[\text{payoff}^*] = \sum_{i=1}^m \alpha_i^* \left(\sum_{j=1}^n \sum_{k=1}^s y_{i,j,k} x_{j,k}^* \right)$. Suppose that we approximate the vector α^* by a vector $\alpha_\varepsilon = (\alpha_{1,\varepsilon}, \dots, \alpha_{m,\varepsilon})$ such that $|\alpha_i^* - \alpha_{i,\varepsilon}| \leq \varepsilon$ for each i . Then, if $\varepsilon \leq \frac{\delta}{nsy}$ we can compute a solution with a total expected payoff of at least $\mathbb{E}[\text{payoff}] - \delta$.

Proof. Our algorithm is simple. Plugging the values of this α_ε in (Q2) reduces it to a linear program, which can be solved optimally in polynomial time giving a solution vector, say \mathbf{x}_ε . Our solution vectors are α_ε and \mathbf{x}_ε . Obviously, all the constraints are satisfied, so we just need to check the total expected payoff of our solution. For notational convenience, let $F(\alpha, \mathbf{x}) = \sum_{i=1}^m \alpha_i \left(\sum_{j=1}^n \sum_{k=1}^s y_{i,j,k} x_{j,k} \right)$ for two vectors $\mathbf{x} = (x_{1,1}, \dots, x_{1,s}, x_{2,1}, \dots, x_{2,s}, \dots, x_{n,1}, \dots, x_{n,s})$ and $\alpha = (\alpha_1, \dots, \alpha_m)$; thus $F(\alpha^*, \mathbf{x}^*) = \mathbb{E}[\text{payoff}^*]$. Then,

$$\begin{aligned} |F(\alpha^*, \mathbf{x}^*) - F(\alpha_\varepsilon, \mathbf{x}^*)| &\leq \varepsilon \sum_{j=1}^n \sum_{k=1}^s y_{i,j,k} \leq \varepsilon nsy \\ \implies F(\alpha_\varepsilon, \mathbf{x}_\varepsilon) &\geq F(\alpha_\varepsilon, \mathbf{x}^*) \geq F(\alpha^*, \mathbf{x}^*) - \varepsilon nsy \geq F(\alpha^*, \mathbf{x}^*) - \delta \end{aligned}$$

□

To get such a α_ε , for every $\alpha_{i,\varepsilon}$ we try out all rational numbers between 0 and 1 of the form $\frac{j\delta}{2nsy}$ for $j = 0, 1, \dots, \frac{2nsy}{\delta}$ until we succeed. The total number of choices is at most $\left(\frac{2nsy}{\delta} + 1 \right)^m$, which is **pseudo-polynomial**¹⁹ in the size of the input since m is fixed.

The result for INT-MULTI-SSBO follows by using the above proof with Lemma 4.

(c) When $ns = O(\log m)$ then we can try out all possible $\text{poly}(m)$ assignments of keywords to slots. For each assignment, we can directly calculate the values of $\alpha_1, \alpha_2, \dots, \alpha_m$. We take the best of all

¹⁹The running time is *not* strongly polynomial since the input size depends polynomial on $\log_2 y$ (see Section 2.5).

such solutions.

(d) Let $p_1(n)$ be a polynomial in n such that $\max \left\{ y, \max_i \{B_i\}, \max_{i,j,k} \{w_{i,j,k}\} \right\} < p_1(n)$. By the proof in part (a), to ensure an absolute error of δ , it suffices to try all vectors $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_m)$ in which each α_i is a non-negative rational number with numerator and denominator at most $p_2(n)$ for some polynomial $p_2(n)$, and provide a solution of INT-MULTI-SSBO for this α in polynomial time. We will refer to B_i as the “expected budget” for the i th scenario. Let $\mathbb{E}[\text{payoff}(j, k, b_1, \dots, b_m)]$ be the optimal value of the expected payoff when *no* slot was selected after the k th slot of the j th keyword and the expected budget for the i th scenario was b_i . It is easy to see that the following recurrence holds:

$$\mathbb{E}[\text{payoff}(j, k, b_1, \dots, b_m)] = \max \left\{ \sum_{i=1}^m y_{i,j,k} + \mathbb{E}[\text{payoff}(j-1, s, b_1 - \alpha_1 w_{1,j,k}, \dots, b_m - \alpha_m w_{m,j,k})], \right. \\ \left. \mathbb{E}[\text{payoff}(j, k-1, b_1, \dots, b_m)] \right\}$$

Based on the above recurrence, it is easy to design a polynomial time dynamic programming algorithm to compute the optimal solution $\mathbb{E}[\text{payoff}(n, s, B_1, \dots, B_m)]$ of INT-MULTI-SSBO. \square

7.2 Limitations of the Semidefinite Programming Relaxation Approaches for SSBO

A natural Semidefinite programming (SDP) relaxation approach to solve quadratic programs such as (Q1), extensively used in existing literatures for efficient approximations of quadratic programs for MAX-CUT, MAX-2SAT and many other problems [24], is as follows. We first add some redundant inequalities to (Q1). For every i and j we add the inequality $\alpha_i x_j \geq 0$. Clearly, this does not change the solutions of (Q1). Then, (Q1) can be relaxed to a vector program (V) by replacing the variables by $(m+n)$ -dimensional vectors and the product of variables by the inner product (denoted by \cdot) of the corresponding vectors. The resulting vector program is shown in Fig. 6; it is well known that (V) is a relaxation of (Q1) (e.g., see [24]).

Since the lower bounds in Theorem 5 have $\varepsilon < 1$ and thus leaves a “very small” gap between this lower bound and the upper bound in Theorem 1, one might wonder if the gap can be somewhat narrowed down by designing an approximation algorithm based on the SDP-relaxation approaches whose approximation ratio is, say, $o\left(\frac{m}{\log m}\right)$ or $o\left(\frac{\log d_n}{\log \log d_n}\right)$? However, we show that the large integrality gap of the SDP-relaxation does not allow for such a possibility.

Since the lower bounds in Theorem 5 have $\varepsilon < 1$ and thus leaves a “very small” gap between this lower bound and the upper bound in Theorem 1, one might wonder if the gap can be somewhat narrowed down by designing an approximation algorithm based on the SDP-relaxation approaches whose approximation ratio is, say, $o\left(\frac{m}{\log m}\right)$ or $o\left(\frac{\log d_n}{\log \log d_n}\right)$? However, we show that the large integrality gap of the SDP-relaxation does not allow for such a possibility.

Lemma 18 (Limitations of SDP-relaxation approaches). *Let $\kappa = 1$. Let OPT_{Q1} and OPT_V be the total optimal payoff for an instance of (Q1) and the optimal value of the objective function of (V), respectively. Then, $\frac{OPT_V}{OPT_{Q1}} \geq \frac{m}{2} = \Theta\left(\frac{\log d_n}{\log \log d_n}\right)$.*

$$\begin{aligned} & \text{(* Vector program (V) *)} \\ & \text{maximize } \sum_{i=1}^m \sum_{j=p}^q y_{i,j} \mathcal{U}_i \cdot \mathcal{V}_j \\ & \text{subject to } \forall 1 \leq i \leq m: \sum_{j=1}^n c_{i,j} y_{i,j} \mathcal{U}_i \cdot \mathcal{V}_j \leq B_i \\ & \quad \forall 1 \leq i \leq m: \forall 1 \leq j \leq n: \mathcal{U}_i \cdot \mathcal{V}_j \geq 0 \\ & \quad \forall 1 \leq i \leq m: \mathcal{U}_i \cdot \mathcal{U}_i \leq 1 \\ & \quad \forall 1 \leq i \leq m: \mathcal{U}_i \in \mathbb{R}^{m+n} \\ & \quad \forall 1 \leq j \leq n: \mathcal{V}_j \cdot \mathcal{V}_j \leq 1 \\ & \quad \forall 1 \leq j \leq n: \mathcal{V}_j \in \mathbb{R}^{m+n} \end{aligned}$$

Figure 6: SDP-relaxation of (Q1).

Proof. We reuse the notations and terminologies used in the proof of Theorem 5. Let the given graph G be a completely connected graph; thus $\Delta_{\text{ind}} = 1$. We construct an instance of SSBO as in Theorem 5. Thus, $\Delta_{Q1} < 1 + \Delta_{\text{ind}} = 2$. Note that $c_n = d_n = m^{6m}$ and thus $m = \Theta(\log d_n / \log \log d_n)$.

However, we show that $\text{OPT}_{\text{vector}} \geq m$. Let $\mathcal{U}_1, \dots, \mathcal{U}_m$ be a set of mutually orthogonal *unit-norm* vectors in \mathbb{R}^{m+n} and let $\mathcal{V}_i = \mathcal{U}_{m-i+1}$ for $1 \leq i \leq m$. Thus, $\mathcal{U}_i \cdot \mathcal{V}_j$ is 1 if $i + j = m + 1$ and is 0 otherwise, and $\mathcal{U}_i \cdot \mathcal{U}_i = \mathcal{V}_i \cdot \mathcal{V}_i = 1$ for all i . Obviously, $\sum_{i=1}^m \sum_{j=1}^n y_{i,j} \mathcal{U}_i \cdot \mathcal{V}_j = m$. We now verify that this is indeed a valid solution of **(V)** by checking that it satisfies all the constraints $\left(\sum_{j=1}^n w_{i,j} \mathcal{U}_i \cdot \mathcal{V}_j\right) \leq B_i$ for $1 \leq i \leq n$. It can be seen that $\left(\sum_{j=1}^n w_{i,j} \mathcal{U}_i \cdot \mathcal{V}_j\right) = w_{i,m-i+1} = c_{m-i+1} = B_i$. \square

7.3 Combinatorial Dual of SSBO Problems

In DUAL-SSBO, the natural combinatorial dual version of SSBO, we are given a lower bound, say P , on $\mathbb{E}[\text{payoff}]$. Our goal is to compute the *minimum* possible value of the budget B of the advertiser such that his/her total expected payoff is *at least* P . The dual version DUAL-MULTI-SSBO of MULTI-SSBO can be defined in a manner analogous to that of DUAL-SSBO. DUAL-SSBO can be reformulated as the quadratic program **(Dual-Q1)** shown in Fig. 7.

(* Quadratic program **(Dual-Q1)** *)

minimize B

subject to $\sum_{i=1}^m \sum_{j=1}^n \alpha_i x_j y_{i,j} \geq P$

$\forall 1 \leq i \leq m: \alpha_i \left(\sum_{j=1}^n w_{i,j} x_j\right) \leq \varepsilon_i B$

$\forall 1 \leq i \leq m: 0 \leq \alpha_i \leq 1$

$\forall 1 \leq j \leq n: 0 \leq x_j \leq 1$

Figure 7: Quadratic program for DUAL-SSBO.

Obviously, DUAL-SSBO is NP-hard since SSBO is NP-hard. For a given required expected profit P , let \mathcal{B}_P be the minimum budget that achieves the expected total profit P . We define a *bi-criteria approximation* for DUAL-SSBO in the following manner:

a (δ, γ) -approximation for DUAL-SSBO, for $\delta, \gamma \geq 1$, is a solution that achieves an expected total profit of at least $\frac{P}{\delta}$ with a budget of $\gamma \mathcal{B}_P$.

Lemma 19.

(a) (Inapproximability of DUAL-SSBO via inapproximability of SSBO)

- If FRAC-SSBO cannot be approximated to within a ratio of $\rho > 1$ for some parameter range, then DUAL-FRAC-SSBO also cannot be approximated to within a ratio of ρ for the same parameter range.
- If INT-SSBO cannot be approximated to within a ratio of $\rho > 1$ for some parameter range, then DUAL-INT-SSBO also cannot be approximated to within a ratio of $\frac{\rho}{200 \ln m}$ for the same parameter range.

(b) (Bi-criterion approximation of DUAL-FRAC-SSBO via FRAC-SSBO) If FRAC-SSBO can be approximated to within a ratio of $\rho > 1$ for some parameter range, then FRAC-SSBO has a $(\rho, 1)$ -approximation in the same parameter range.

Proof. Let $\mathbb{E}[\text{payoff}^{\mathcal{B}}]$ be the optimal total expected payoff for SSBO when the budget is \mathcal{B} . For any constant $\Delta > 1$, a solution of **(Q1)** with a budget of \mathcal{B} is obviously also a solution of the same instance of **(Q1)** with a budget of $\Delta \mathcal{B}$. This implies $\mathbb{E}[\text{payoff}^{\Delta \mathcal{B}}] \geq \mathbb{E}[\text{payoff}^{\mathcal{B}}]$. Let

$p = \sum_{i=1}^m \sum_{j=1}^n y_{i,j}$ and $b = \max_{1 \leq i \leq m} \left\{ \sum_{j=1}^n a_{i,j} c_{i,j} \right\}$; note that both $\log_2 p$ and $\log_2 b$ are polynomial in the size of the input (see Section 2.5).

We prove (a) by contradiction. Suppose that some version of DUAL-SSBO has a ρ -approximation. Consider an instance of the *same* version of SSBO and suppose the budget is B . We do a binary search in the range of positive integers $[1, p]$ in polynomial time with the approximation algorithm for DUAL-SSBO to find a $\mathcal{P} \in [1, p]$ such that $\mathcal{B}_{\mathcal{P}-1} < \rho B$ but $\mathcal{B}_{\mathcal{P}} \geq \rho B$. Consider this solution of DUAL-SSBO and suppose that \mathcal{B}^* is the actual optimal value of the budget corresponding to the total expected payoff \mathcal{P} . Thus, $\mathcal{B}^* \geq \frac{\mathcal{B}_{\mathcal{P}}}{\rho} \geq B$ and $\mathbb{E}[\text{payoff}^{\mathcal{B}_{\mathcal{P}}}] \geq \mathbb{E}[\text{payoff}^{\mathcal{B}^*}] \geq \mathbb{E}[\text{payoff}^B]$. Suppose that we now divide every x_i by ρ . This provides a valid solution of FRAC-SSBO with a total expected payoff of at least $\frac{\mathbb{E}[\text{payoff}^{\mathcal{B}_{\mathcal{P}}}]}{\rho}$. By Lemma 4, from this valid solution of FRAC-SSBO one can obtain a solution of INT-SSBO with a total expected payoff of at least $\frac{\mathbb{E}[\text{payoff}^{\mathcal{B}_{\mathcal{P}}}]}{200 \rho \ln m}$.

To prove (b), suppose that some version of SSBO with a budget of \mathcal{B} has a ρ -approximation algorithm. Consider an instance of the same version of DUAL-SSBO with a requirement of total expected payoff of \mathcal{P} and let $\mathcal{B}_{\mathcal{P}}$ be the value of an optimal budget for this instance. Since $\left(1 - \frac{1}{\mathcal{B} + 1}\right) \mathbb{E}[\text{payoff}^{\mathcal{B}+1}] \leq \mathbb{E}[\text{payoff}^{\mathcal{B}}] \leq \mathbb{E}[\text{payoff}^{\mathcal{B}+1}]$, we do a binary search in the range of positive integers $[1, b]$ in polynomial time with the ρ -approximation algorithm for SSBO to find a $\mathcal{B} \in [1, b]$ such that $\frac{\mathcal{P}}{\rho} \leq \mathbb{E}[\text{payoff}^{\mathcal{B}}] \leq \rho \mathcal{P} + 1$. Thus, this provides a solution of the DUAL-SSBO with a total expected payoff of at least $\frac{\mathcal{P}}{\rho}$ and a budget of at most $\mathcal{B}_{\mathcal{P}}$, giving the desired $(\rho, 1)$ -approximation in (b). \square

8 Conclusion

We have presented the first known approximation algorithms as well as hardness results for stochastic budget optimization under the scenario model. The scenario model is natural in many areas, and it is particularly apt for internet ad systems. We obtained our results by making the connection between these problems and a special case of bipartite quadratic programs; we exploited this intuition crucially in both approximation algorithms and hardness proofs. These class of quadratic programs may have independent applications elsewhere.

Our work shows that there are several instances of parameters where stochastic budget optimizations are solvable with reasonable computational resource even with multiple slots. Our hope is that therefore, in practice, one can carefully model particular applications such as sponsored search, so that the parameters are suitable, and advertisers can optimize their campaigns more effectively than is typically done now by applying some of the algorithms in this paper.

Acknowledgements

We thank the reviewers for their detailed comments which improved both the readability and the technical content of the paper.

References

- [1] N. Alon and A. Naor. *Approximating the cut-norm via Grothendieck's inequality*, 36th ACM Symposium on Theory of Computing, 72-80, 2004.

- [2] G. Aggarwal, A. Goel, and R. Motwani. *Truthful auctions for pricing search keywords*, 7th ACM Conference on Electronic Commerce, 1-7, 2006.
- [3] C. Borgs, J. Chayes, O. Etesami, N. Immorlica, K. Jain and M. Mahdian. *Dynamics of bid optimization in online advertisement auctions*, 16th International World Wide Web Conference, 531-540, 2007.
- [4] T.-M. Bu, X. Deng and Q. Qi. *Forward looking Nash equilibrium for keyword auction*, Information Processing Letters, 105(2), 41-46, 2008.
- [5] M. Cary, A. Das, B. Edelman, I. Giotis, K. Heimerl, A. R. Karlin, C. Mathieu and M. Schwarz. *Greedy bidding strategies for keyword auctions*, 8th ACM conference on Electronic commerce, 262-271, 2007.
- [6] C. Chekuri and S. Khanna. *A PTAS for the multiple knapsack problem*, 11th ACM-SIAM Symposium on Discrete Algorithms, 213-222, 2000.
- [7] F. Cocco, A. Consiglio and S. Zenios. *Scenario Optimization Asset and Liability Modeling for Endowments with Guarantees*, Wharton Financial Institutions Center WP No. 00-41, June 12, 2000, available at <http://ssrn.com/abstract=253039> or doi:10.2139/ssrn.253039.
- [8] R. S. Dembo. *Scenario optimization*, Annals of Operations Research, 30 (1), 63-80, 1991.
- [9] B. Edelman, M. Ostrovsky, and M. Schwarz. *Internet advertising and the generalized second price auction: Selling billions of dollars worth of keywords*, American Economic Review, 97(1), 242-259, 2007.
- [10] U. Feige. *A threshold of for approximating set cover*, Journal of the ACM, 45, 634-652, 1998.
- [11] U. Feige and L. Lovasz. *Two-prover one-round proof systems: their power and their problems*, 24th ACM Symposium on Theory of Computing, 733-744, 1992.
- [12] J. Feldman, S. Muthukrishnan, M. Pál and C. Stein. *Budget optimization in search-based advertising auctions*, 8th ACM Conference on Electronic Commerce, 40-49, 2007.
- [13] M. R. Garey and D. S. Johnson. *Computers and Intractability - A Guide to the Theory of NP-Completeness*, W. H. Freeman & Co., 1979.
- [14] M.X. Goemans and D.P. Williamson. *Improved Approximation Algorithms for Maximum Cut and Satisfiability Problems Using Semidefinite Programming*, Journal of the ACM, 42, 1115-1145, 1995.
- [15] J. Håstad. *Clique is hard to approximate within $n^{1-\epsilon}$* , Acta Mathematica, 182, 105-142, 1999.
- [16] H. Mausser and D. Rosen. *Applying Scenario Optimization to Portfolio Credit Risk*, The Journal of Risk Finance, 2 (2), 36-48, 2001.
- [17] K. Jansen. *Parameterized approximation scheme for the multiple knapsack problem*, 20th ACM-SIAM Symposium on Discrete Algorithms, 665-674, 2009.
- [18] G. Kortsarz. *On the hardness of approximating spanners*, Algorithmica, 30 (3), 432-450, 2001.
- [19] E. L. Lawler. *Fast Approximation Algorithms for Knapsack Problems*, Mathematics of Operations Research, 4(4), 339-356, 1979.

- [20] R. Motwani and P. Raghavan. *Randomized Algorithms*, Cambridge University Press, 1995.
- [21] S. Muthukrishnan, M. Pál and Z. Svitkina. *Stochastic Models for Budget Optimization in Search-Based Advertising*, Algorithmica, 58 (4), 1022-1044, 2010.
- [22] R. Raz. *A parallel repetition theorem*, SIAM Journal of Computing, 27 (3), 763-803, 1998.
- [23] H. Varian. *Position auctions*, International Journal of Industrial Organization, 25(6), 1163-1178, 2007.
- [24] V. Vazirani. *Approximation Algorithms*, Springer-Verlag, 2001.
- [25] S. A. Zenios (editor). *Financial Optimization*, Cambridge University Press, 1996.
- [26] Y. Zhou, D. Chakrabarty, R. Lukose. *Budget Constrained Bidding in Keyword Auctions and Online Knapsack Problems*, 4th International Workshop On Internet And Network Economics, 566-576, 2008.
- [27] Y. Zhou and V. Naroditskiy. *Algorithm for Stochastic Multiple-Choice Knapsack Problem and Application to Keywords Bidding*, 17th International World Wide Web Conference, poster paper.