

Random Knapsack in Expected Polynomial Time

Rene Beier*
rbeier@mpi-sb.mpg.de
Max-Planck-Institut für Informatik
Saarbrücken, Germany

Berthold Vöcking†
voecking@ls2.cs.uni-dortmund.de
Fachbereich Informatik
Universität Dortmund, Germany

ABSTRACT

In this paper, we present the first average-case analysis proving an expected polynomial running time for an exact algorithm for the 0/1 knapsack problem. In particular, we prove, for various input distributions, that the number of *dominating solutions* (i.e., Pareto-optimal knapsack fillings) to this problem is polynomially bounded in the number of available items. An algorithm by Nemhauser and Ullmann can enumerate these solutions very efficiently so that a polynomial upper bound on the number of dominating solutions implies an algorithm with expected polynomial running time.

The random input model underlying our analysis is very general and not restricted to a particular input distribution. We assume adversarial weights and randomly drawn profits (or vice versa). Our analysis covers general probability distributions with finite mean, and, in its most general form, can even handle different probability distributions for the profits of different items. This feature enables us to study the effects of correlations between profits and weights. Our analysis confirms and explains practical studies showing that so-called *strongly correlated* instances are harder to solve than *weakly correlated* ones.

Categories and Subject Descriptors

G.1.6 [Optimization]: Integer Programming; F.2.0 [Analysis of Algorithms and Problem Complexity]: Miscellaneous

General Terms

Algorithms

Keywords

knapsack problem, exact algorithms, average case analysis

*Supported by the graduated studies program on “Quality Guarantees for Computer Systems” funded by the German Science Foundation (DFG).

†Supported in part by DFG grant Vo889/1-1.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

STOC’03, June 9–11, 2003, San Diego, California, USA.
Copyright 2003 ACM 1-58113-674-9/03/0006 ...\$5.00.

1. INTRODUCTION

The 0/1 knapsack problem is one of the most intensively studied combinatorial optimization problems having a wide range of applications in industry and financial management, e.g., cargo loading, cutting stock, budget control. The problem is defined as follows. Given n items with positive weights w_1, \dots, w_n and profits p_1, \dots, p_n and a knapsack capacity c , find a subset $S \subseteq [n] := \{1, 2, \dots, n\}$ such that $\sum_{i \in S} w_i \leq c$ and $\sum_{i \in S} p_i$ is maximized. This problem is of special interest not only from a practical point of view but also for theoretical reasons as it can be seen as the simplest possible 0/1 linear program because the set of feasible solutions is described by a single constraint only. Starting with the pioneering work of Dantzig [4], the problem has been studied extensively in practice as well as in theory. In this paper, we are concerned with the average-case analysis of exact algorithms for this problem. The major motivation for our study is the gap between the theoretically proven high worst-case complexity and the observed efficiency of some algorithms on various practical instances.

The knapsack problem is one of those optimization problems for which NP-hardness theory concludes that it is hard to solve in the worst case. Despite the exponential worst-case running times of all known knapsack algorithms, several large scale instances can be solved to optimality very efficiently [1, 7, 10, 11, 19]. In particular, randomly generated instances seem to be quite easy to solve. In order to explain the observed efficiency on random instances, several theoretical studies investigate the structure and complexity of random knapsack instances [12, 14, 15, 2, 6].¹ Although significant progress has been made in understanding the structure and complexity of random knapsack instances, until now there has been no proof that the knapsack problem can be solved in expected polynomial time under any reasonable stochastic input model.

The best known result on the running time of an exact algorithm for the knapsack problem was shown by Goldberg and Marchetti-Spaccamela [12].² They investigate a so-called *core algorithm*, an algorithmic concept suggested by Balas and Zemel [1]. The idea is to start with the optimal fractional solution containing at most one fractional item and then to exchange some of the items until an optimal integral solution is found. The set of items that are candidates to be exchanged, is called the *core set*, and the hope is that

¹Related work deals with the hardness of random instances for knapsack cryptosystems [5, 8, 13]. These cryptosystems, however, are based on the hardness of random subset sum which is not directly affected by our results.

²This work was later extended to the multidimensional knapsack problem by Dyer and Frieze [6].

the size of the core set is relatively small for “typical instances”. In their average-case analysis, Goldberg and Marchetti-Spaccamela study the size and structure of the core set assuming that profits and weights of all items are drawn independently and uniformly at random from the interval $[0, 1]$. As a first step towards analyzing core algorithms, Lueker proved an upper bound on the expected gap between the optimal integral and the optimal fractional solution [14]. Based on this result, Goldberg and Marchetti-Spaccamela were able to prove structural properties of the core set resulting in the following bound on the running time of a Las Vegas type core algorithm. For every fixed $k > 0$, with probability at least $1 - 1/k$, the running time of the algorithm does not exceed a specified upper bound that is polynomial in the number of items. However, the degree of this polynomial is quite large, the leading constant in the exponent is at least a three digit number, and, more dramatically, the degree grows with the reciprocal of the failure probability like $\sqrt{k} \log(k)$. Unfortunately, such a result does not allow to conclude a sub-exponential upper bound on the expected running time of the algorithm.

In this paper, we present the first average-case analysis proving expected polynomial running time for an exact algorithm for the 0/1 knapsack problem. We improve the result of Goldberg and Marchetti-Spaccamela in several other aspects as well. In particular, our random input model is not restricted to the uniform distribution. We assume that the weights of the items are chosen by an adversary and their profits are chosen according to arbitrary probability distributions with finite mean (or vice versa). Moreover, the profits of different items can follow different distributions. This enables us to study the effects of correlations between profits and weights, which is also a major aspect in the more recent practical studies [10, 11, 19]. Furthermore, our analysis does not involve large constants. The degree of the polynomials in our upper bounds on the running time ranges from 3 to 5 depending on the underlying probability distribution.

Our analysis is based on an elegant algorithm presented by Nemhauser and Ullmann [17] in 1969. This algorithm can be viewed as a sparse dynamic programming approach. Its efficiency on practical instances was already mentioned by Nemhauser and Ullmann in their seminal paper but until now this efficiency has been shown only within experimental studies, see e.g. [11]. Before we turn to a more detailed presentation and discussion of our results and techniques, let us introduce this algorithm.

1.1 The Nemhauser/Ullmann algorithm

A brute force method to solve the knapsack problem is to enumerate all possible subsets over the n items. In order to reduce the search space, a domination concept is used which is usually attributed to Weingartner and Ness [22]. A subset $S \subseteq [n]$ with weight $w(S) = \sum_{i \in S} w_i$ and profit $p(S) = \sum_{i \in S} p_i$ dominates another subset $T \subseteq [n]$ if $w(S) \leq w(T)$ and $p(S) \geq p(T)$. For simplicity assume that no two subsets have the same profit. Then no subset dominated by another subset can be an optimal solution to the knapsack problem, regardless of the specified knapsack capacity. Consequently, it suffices to consider those sets that are not dominated by any other set, the so-called *dominating sets*. In other terminology, dominating sets are *Pareto-optimal solutions*, i.e., solutions that cannot be improved in weight and profit simultaneously by other solutions.

Nemhauser and Ullmann [17] introduce the following elegant algorithm computing a list of all dominating sets in an iterative manner. For $i \in [n]$, let $S(i)$ be the sequence of dominating subsets over the items $1, \dots, i$. The sets in $S(i)$ are assumed to be listed in increasing order of their weights. Given $S(i)$, the sequence $S(i+1)$ can be computed from $S(i)$ as follows: First duplicate all subsets in

$S(i)$ and then add item $i+1$ to each of the duplicated sets. In this way we obtain two ordered sequences of sets. Now we merge the two sequences by removing the sets dominated by any other set in the union of the two sequences. The result is the ordered sequence $S(i+1)$ of dominating sets over the items $1, \dots, i+1$.

For the purpose of illustration and a better understanding, let us take a different view on this algorithm. For $i \in [n]$, let $f_i : \mathbb{R} \rightarrow \mathbb{R}$ be a mapping from weights to profits such that $f_i(t)$ is the maximum profit over all subsets of $[i]$ with weight at most t . Observe that f_i is a non-decreasing step function changing its value at those weights that correspond to dominating subsets. In particular, the number of steps in f_i equals the number of dominating sets over the items in $[i]$. Figure 1 (a) shows such a step function for a small instance, (b) shows a step function f_i and the copy of this step function which is implicitly generated in the algorithm described above, finally, (c) illustrates how the two step functions are merged: the graph of the resulting step function f_{i+1} is simply the upper envelope of the graph f_i and its shifted copy.

For the model of computation, let us assume a uniform RAM that can add and compare numbers in constant time. Then the sequence S_{i+1} can be calculated from the sequence S_i in time linear in the length of S_i , that is, linear in the number of dominating subsets over the items $1, \dots, i$. Since the optimal knapsack filling is described by one of the subsets in the list S_n , namely the subset with the largest weight not exceeding the capacity, generating S_n basically solves the knapsack problem. This yields the following lemma:

LEMMA 1. *For every $i \in [n]$, let $q(i)$ denote an upper bound on the (expected) number of dominating sets over the items in $1, \dots, i$, and assume $q(i+1) \geq q(i)$. The Nemhauser/Ullmann algorithm computes an optimal knapsack filling in (expected) time*

$$O(\sum_{i=1}^n q(i)) = O(n \cdot q(n)).$$

In the worst case, the number of dominating sets is 2^n . This case occurs when profits and weights are identical. In fact, this instance is also the worst case for the core algorithm of Goldberg and Marchetti-Spaccamela. In general, the running time of the Nemhauser/Ullmann algorithm can be bounded from above by $O(\sum_{i=1}^n 2^i) = O(2^{n+1})$. However, if weights and profits are independent or only weakly correlated, experiments show that the number of dominating sets and, hence, the running time, is much smaller [11]. Furthermore, if the input numbers are positive integers, then the running time of the algorithm can be bounded pseudo-polynomially as, in this case, step function f_i can have at most iP steps, with P denoting the maximum profit of an individual item. At this point, let us remark that the dominating set algorithm can also be seen as a sparse dynamic programming approach.

Horowitz and Sahni [18] present a nice variation of the algorithm above achieving a better worst-case performance by splitting the set of items into two groups and processing each of them separately. In this way, they can achieve a square root improvement on the worst-case running time, that is, the worst-case running time is reduced to $O(2^{n/2})$. This improvement, however, does not translate to the average case and, therefore, we will not further follow this approach but focus on the original algorithm by Nemhauser and Ullmann. The challenge in the analysis of this algorithm is to estimate the number of dominating sets over a set of items with randomly drawn profits and to determine how the correlation between profits and weights influences the number of dominating sets.

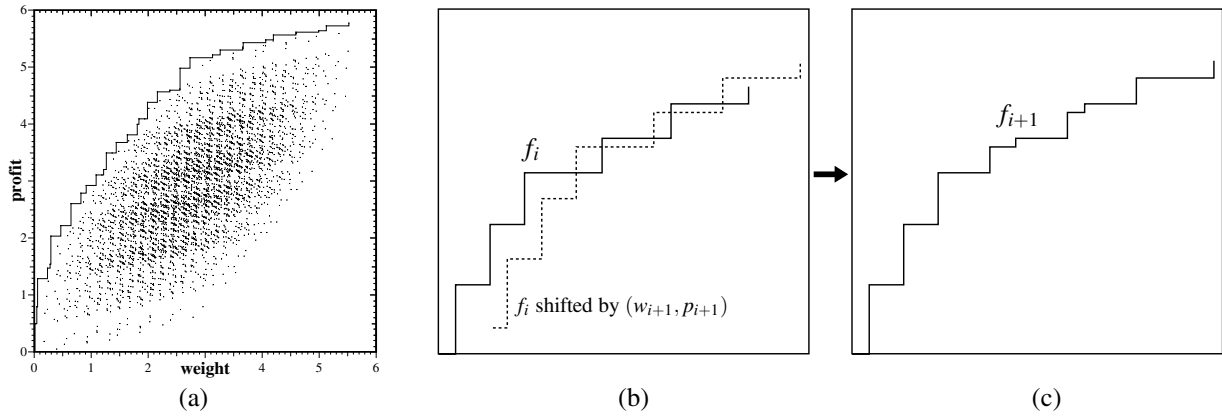


Figure 1: (a) Profits and weights of all subsets over 12 items and corresponding step function f . Weight and profit of each item are independent random numbers from $[0, 1]$. (b) Shift step function f_i by vector (w_{i+1}, p_{i+1}) . (c) Upper envelope gives step function f_{i+1} .

1.2 Our contribution

In this paper, we present the first average-case analysis proving expected polynomial running time for an exact algorithm solving the 0/1 knapsack problem. Our analysis is very robust and can be applied to several random input models. In general, we assume that the weights of the items and the capacity of the knapsack are specified by an adversary. The profits are chosen according to various classes of continuous as well as discrete probability distributions.³

In particular, we consider the following classes of continuous probability distributions:

- *The uniform distribution:* Suppose profits are chosen independently, uniformly at random from the interval $[0, 1]$. Let q denote the number of dominating sets over n items. We prove $\mathbf{E}[q] = O(n^3)$. Combining this bound with Lemma 1 immediately implies an upper bound of $O(n^4)$ on the expected running time of the Nemhauser/Ullmann algorithm. More details are explained in Section 2.
- *The exponential and other long tailed distributions:* If profits are chosen according to the exponential distribution, then $\mathbf{E}[q] = O(n^2)$ and, hence, $\mathbf{E}[T] = O(n^3)$. Furthermore, we can generalize this upper bound towards all continuous distributions with finite mean whose tails (under appropriate normalization) are not shorter than the tail of the exponential distribution. For example, for the Pareto distribution with parameter $a > 1$, $\mathbf{E}[q] = O(\frac{a}{a-1} n^2)$ and, hence, $\mathbf{E}[T] = O(\frac{a}{a-1} n^3)$. Moreover, our analysis allows even a heterogeneous mix of long-tailed distributions, that is, different distributions for the profits of different items. For more details about these results, see Section 3.
- *Distributions with non-increasing density:* For these distributions we can show a lower bound. Suppose profits are chosen according to an arbitrary probability distribution with non-increasing density function. We prove that there is a way to choose the weights such that $\mathbf{E}[q] = \Omega(n^2)$. Thus our upper bound for the exponential distribution is tight. Further details can be found in Section 4.

³Symmetrically, our analysis can be applied to adversarial profits and random weights, too. Details will be explained in a full version of this paper.

- *General, continuous distributions:* Suppose profits are chosen according to arbitrary, possibly different probability distributions with finite mean. Let μ denote the maximum expected profit over all items and ϕ the *maximum density*, i.e., the maximum value taken by any of the density functions describing the probability distributions for the profits. Then $\mathbf{E}[q] = O(\mu\phi n^4)$ and $\mathbf{E}[T] = O(\mu\phi n^5)$. This very general result is presented in Section 5.

The result for general distributions allows to study the influence of the degree of randomness in the problem specification on the complexity of the knapsack problem. Let us normalize the specified profit distributions by multiplying all profits with $\frac{1}{\mu}$. In this way the maximum expected profit is 1 and the upper bound on the expected number of dominating sets simplifies to $\mathbf{E}[q] = O(\phi n^4)$. Under this normalization, the maximum density ϕ can be seen as a parameter describing how much randomness is available. For $\phi \rightarrow \infty$ the randomness in the specification of the profits goes to zero and, in this case, an adversary can specify an input such that the expected number of dominating sets is exponential. In contrast, if ϕ is bounded from above by some constant term (as, e.g., in the case of all distributions mentioned above), then the described instances inhabit a high degree of randomness and the expected number of dominating sets is polynomial.

Furthermore, this result enables us to study the effects of correlations between profits and weights. (This aspect has also been considered in several recent practical studies [10, 11, 19].) For this purpose, let us present our result for general distributions in form of a so-called “smoothed analysis”. Spielman and Teng [21] introduced the following random input model that allows to perform a mix of worst-case and average-case analysis. The idea is to initially start with a worst-case instance and then to introduce randomization by perturbing the profits according to the Gaussian distribution with small standard deviation. More precisely, initially all profits and weights are specified by an adversary such that all these numbers fall into the interval $[0, 1]$. Then profits are perturbed using the Gaussian distribution with some specified standard deviation σ . Of course, for our application, the random perturbation might produce a small number of negative profits. These profits, however, can be removed immediately as they obviously do not belong to the optimal knapsack filling. The resulting distributions for the profits of the surviving items have mean $\mu \leq 1 + \sigma$ and maximum density $\phi \leq 1/\sqrt{2\pi\sigma^2} < 1/\sigma$. For such a perturbed instance, our

analysis immediately implies $\mathbf{E}[q] = O(n^4/\sigma + n^4)$. We observe that if the adversary chooses profits equal to weights, then the unperturbed instance is completely correlated and has 2^n dominating sets. Now perturbing the profits decreases the correlation. In particular, a small correlation corresponds to a large standard deviation σ and this in turn implies a small upper bound on the expected number of dominating sets. In other words, the complexity of the problem diminishes when decreasing the correlation between profits and weights.

Finally, we want to point out that our results can also be generalized towards discrete probability distributions. Interestingly, this gives a complete tradeoff ranging from pseudo-polynomial running time for worst-case inputs without randomness to polynomial running time for fully random instances. For more details refer to Section 6.

2. THE UNIFORM DISTRIBUTION

In this section, profits are assumed to be chosen uniformly at random from $[0, 1]$ and, as in all of our analyses, weights are chosen by an adversary. The following upper bound on the expected number of dominating sets combined with the result in Lemma 1 implies an upper bound of $O(n^4)$ on the expected running time of the Nemhauser/Ullmann algorithm.

THEOREM 2. *Suppose the weights are arbitrary positive numbers and profits are chosen according to the uniform distribution over $[0, 1]$. Let q denote the number of dominating sets over all n items. Then $\mathbf{E}[q] = O(n^3)$.*

PROOF. Let $m = 2^n$ and let S_1, \dots, S_m denote the sequence of all subsets of $[n]$ listed in non-decreasing order of their weights. Let the profit of subset S_u be $P_u = \sum_{i \in S_u} p_i$. For any $2 \leq u \leq m$, define $\Delta_u = \max_{v \in [u]} P_v - \max_{v \in [u-1]} P_v \geq 0$. Observe that S_1 is always a dominating set. For all $u \geq 2$, S_u is dominating if and only if $\Delta_u > 0$. The following lemma shows that the expected increase in profit at dominating sets is $\Omega(n^{-2})$. In other words, the expected height of the steps in the step function f_n is quite large.

LEMMA 3. *For every $u \in \{2, \dots, m\}$, $\mathbf{E}[\Delta_u | \Delta_u > 0] \geq \frac{1}{32n^2}$.*

PROOF. Fix $u \in \{2, \dots, m\}$. Observe that

$$\mathbf{E}[\Delta_u | \Delta_u > 0] \geq \Pr \left[\Delta_u \geq \frac{1}{16n^2} \mid \Delta_u > 0 \right] \cdot \frac{1}{16n^2}.$$

Hence, it suffices to show $\Pr \left[\Delta_u \geq \frac{1}{16n^2} \mid \Delta_u > 0 \right] \geq \frac{1}{2}$. For every $v \in [u-1]$, define $X_v = S_u \setminus S_v$ and $Y_v = S_v \setminus S_u$. It holds

$$\begin{aligned} & \Pr \left[\Delta_u \geq \frac{1}{16n^2} \mid \Delta_u > 0 \right] \\ &= \Pr \left[\forall v : \sum_{i \in S_u} p_i \geq \sum_{i \in S_v} p_i + \frac{1}{16n^2} \mid \forall v : \sum_{i \in S_u} p_i > \sum_{i \in S_v} p_i \right] \\ &= \Pr \left[\forall v : \sum_{i \in X_v} p_i \geq \sum_{i \in Y_v} p_i + \frac{1}{16n^2} \mid \forall v : \sum_{i \in X_v} p_i \geq \sum_{i \in Y_v} p_i \right], \end{aligned} \quad (1)$$

where the universal quantifier ranges over all elements $v \in [u-1]$. Since we consider continuous probability distributions, the relaxation of the strict inequality in the conditioning part does not effect the probability. W.l.o.g., $S_u = [k]$. We distinguish two classes of random variables, namely $\{p_1, \dots, p_k\}$ and $\{p_{k+1}, \dots, p_n\}$. Observe that the X_v 's are subsets of the first class and the Y_v 's are subsets of the second class. For a moment, let us assume that the variables in the second class are fixed arbitrarily. We investigate

the variables in the first class under this assumption. Regardless of how the variables in the second class are fixed, it is unlikely that one of the variables in the first class is much smaller than n^{-1} . In particular,

$$\begin{aligned} & \Pr \left[\exists j \in [k] : p_j \leq \frac{1}{4n} \mid \forall v : \sum_{i \in X_v} p_i \geq \sum_{i \in Y_v} p_i \right] \\ & \leq \sum_{j \in [k]} \Pr \left[p_j \leq \frac{1}{4n} \mid \forall v : \sum_{i \in X_v} p_i \geq \sum_{i \in Y_v} p_i \right] \\ & \leq \sum_{j \in [k]} \Pr \left[p_j \leq \frac{1}{4n} \right] = \frac{k}{4n} \leq \frac{1}{4}. \end{aligned}$$

From now on, we assume $p_j \geq \frac{1}{4n}$, for every $j \in [k]$. Let $L_v = \sum_{i \in X_v} p_i$. Under our assumption, $L_v \geq \frac{1}{4n}$, for every $v \in [u-1]$, because each set X_v contains at least one element of size at least $\frac{1}{4n}$. In the following we will assume that the L_v 's are fixed in a way satisfying this property but, otherwise, arbitrarily. (We will not consider the variables p_1, \dots, p_k anymore.) Under our assumption on the L_v 's, we analyze $\Pr \left[\Delta_u < \frac{1}{16n^2} \mid \Delta_u > 0 \right]$. In order to compensate for the case when our assumption fails, we prove that this probability is at most $\frac{1}{4}$ instead of $\frac{1}{2}$. In particular, using the definition of the L_v 's to rewrite Equation 1, we prove

$$\Pr \left[\forall v : \sum_{i \in Y_v} p_i \leq L_v - \frac{1}{16n^2} \mid \forall v : \sum_{i \in Y_v} p_i \leq L_v \right] \geq \frac{3}{4},$$

for arbitrarily fixed $L_v \geq \frac{1}{4n}$, where the probability refers solely to the random choices for the variables p_{k+1}, \dots, p_n .

Let us now switch to a geometric interpretation. Consider the $(n-k)$ -dimensional polytopes.

$$\begin{aligned} A &= \left\{ (p_{k+1} \times \dots \times p_n) \in [0, 1]^{n-k} \mid \forall v : \sum_{i \in Y_v} p_i \leq L_v - \frac{1}{16n^2} \right\}, \\ B &= \left\{ (p_{k+1} \times \dots \times p_n) \in [0, 1]^{n-k} \mid \forall v : \sum_{i \in Y_v} p_i \leq L_v \right\}. \end{aligned}$$

Figure 2 illustrates the definition of these two polytopes. Clearly, $A \subseteq B$. As we investigate the uniform distribution,

$$\Pr \left[\Delta_u \geq \frac{1}{16n^2} \mid \Delta_u > 0 \right] = \frac{\text{vol}(A \cap B)}{\text{vol}(B)} = \frac{\text{vol}(A)}{\text{vol}(B)},$$

with $\text{vol}(\cdot)$ specifying the volume of the corresponding polytopes. In terms of these volumes, we have to show $\text{vol}(A) \geq \frac{3}{4} \text{vol}(B)$.

At first view, it might seem that the ratio $\text{vol}(A)/\text{vol}(B)$ depends on the number of facets of these polytopes. This number, however, can be exponential, since facets correspond to subsets of $[n]$. Fortunately, however, the following argument shows that the ratio between the two volumes can be estimated in terms of the number of dimensions rather than the number of facets. The idea is to shrink the polytope B uniformly over all dimensions until the shrunken polytope is contained in polytope A . For $\varepsilon \in [0, 1]$, we define

$$B_\varepsilon = \left\{ (p_{k+1} \times \dots \times p_n) \in [0, 1-\varepsilon]^{(n-k)} \mid \forall v : \sum_{i \in Y_v} p_i \leq (1-\varepsilon)L_v \right\}.$$

Obviously, $B = B_0$ and, in general, B_ε can be obtained by shrinking B in each dimension by a factor of $1-\varepsilon$. As the number of dimensions is $n-k$, it holds $\text{vol}(B_\varepsilon) = (1-\varepsilon)^{n-k} \text{vol}(B)$.

Next we investigate how to choose ε such that B_ε is contained in A . Observe that $L_v(1-\frac{1}{4n}) \leq L_v - \frac{1}{16n^2}$ because $L_v \geq \frac{1}{4n}$. Thus

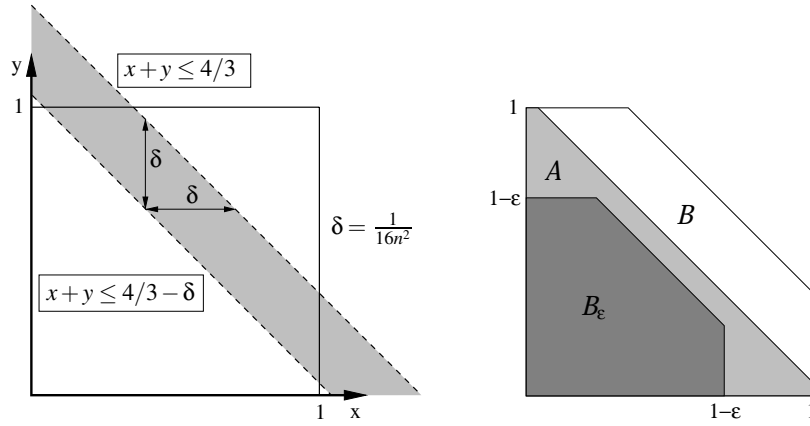


Figure 2: Example for 2-dimensional polytopes A, B and B_ϵ

setting $\epsilon = \frac{1}{4n}$ implies $B_\epsilon \subseteq A$. As a consequence,

$$\begin{aligned} \text{vol}(A) &\geq \text{vol}(B_\epsilon) = (1-\epsilon)^{(n-k)} \text{vol}(B) \\ &\geq (1-\epsilon(n-k)) \text{vol}(B) \geq \frac{3}{4} \text{vol}(B), \end{aligned}$$

which completes the proof of Lemma 3. \square

The lemma above shows that at every dominating set the profit increases by at least $\frac{1}{32n^2}$ on expectation. On the other hand, the expected profit of the knapsack containing all items is $n/2$ as each individual item has profit $1/2$ on expectation. It might be intuitively clear that this implies that the expected number of dominating sets is at most $16n^3$. The following calculation proves this statement in a formal way. Recall that $P_m = \sum_{i \in [n]} p_i$ and $P_1 = \sum_{i \in S_1} p_i = 0$ since $S_1 = \emptyset$. On the one hand,

$$\begin{aligned} \mathbf{E}[P_m] &= P_1 + \sum_{u=2}^m \mathbf{E}[\Delta_u] = \sum_{u=2}^m \Pr[\Delta_u > 0] \cdot \mathbf{E}[\Delta_u \mid \Delta_u > 0] \\ &\geq \sum_{u=2}^m \Pr[\Delta_u > 0] \cdot \frac{1}{32n^2}. \end{aligned}$$

On the other hand, $\mathbf{E}[P_m] = n/2$. Consequently,

$$\mathbf{E}[q] = 1 + \sum_{u=2}^m \Pr[\Delta_u > 0] \leq 1 + 32n^2 \mathbf{E}[P_m] \leq 16n^3 + 1.$$

The additional 1 is due to the empty set S_1 , which is always a dominating set. Thus Theorem 2 is shown.

3. LONG-TAILED DISTRIBUTIONS

One can classify continuous probability distributions by comparing their tails with the tail of the exponential distribution. In principle, if the tail function of a distribution can be lower-bounded by the tail function of the exponential function, then we say the distribution has a “long tail”, and if the tail function can be upper-bounded by the exponential tail function, then we talk about “short tails”. In this section, we investigate the expected number of dominating sets under long-tailed profit distributions. In fact, we can prove a slightly better bound for these distributions than for the short-tailed uniform distribution. Moreover, our analysis can handle heterogeneous distributions. We want to point out that the results we prove for the long-tailed distributions are important tools in the subsequent analysis for general probability distributions.

We need to define the term “long-tailed distribution” more formally. Of special interest for us is the behavior of the tail function under a logarithmic scale. Given any continuous probability distribution with density function $f : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}$, the tail function $T : \mathbb{R}_{\geq 0} \rightarrow [0, 1]$ is defined by $T(t) = \int_t^\infty f(x) dx$. We define the *slope of T at $x \in \mathbb{R}_{\geq 0}$* to be the first derivative of the function $-\ln(T(\cdot))$ at x , i.e., $\text{slope}_T(x) = -[\ln(T(x))]'$. For example, the tail function of the exponential distribution with parameter λ is $T(x) = \exp(-\lambda x)$ so that the slope of this function is $\text{slope}_T(x) = \lambda$, for every $x \geq 0$. In general, $\text{slope}_T(x)$ is a not necessarily continuous function with non-negative real values. The tail of a continuous probability distribution is defined to be *long* if there exists $\alpha > 0$ such that $\text{slope}_T(x) \leq \alpha$, for every $x \in \mathbb{R}_{\geq 0}$.

According to this definition, the exponential distribution has long tails. However, the uniform distribution over $[0, 1]$ (or any other interval) does not have long tails because $\text{slope}_T(x) = 1/(1-x)$, which grows to ∞ for $x \rightarrow 1$. Observe that any distribution with a bounded domain cannot have long tails. A typical example for a distribution with long tails is the Pareto distribution. The tail function of the Pareto distribution with parameter $a > 0$ is $T(x) = (1+x)^{-a}$. These tails are long because $\text{slope}_T(x) = [a \ln(1+x)]' = \frac{a}{1+x} \leq a$, for every $x \geq 0$.

3.1 Analysis

We assume that profits are chosen independently at random according to possibly different long-tailed distributions with finite mean. Weights must be positive but apart from that can be chosen arbitrarily.

THEOREM 4. *For $i \in [n]$, let profit p_i be a random variable with tail function $T_i : \mathbb{R}_{\geq 0} \rightarrow [0, 1]$. Define $\mu_i = \mathbf{E}[p_i]$ and let α_i be an appropriate positive real number satisfying $\text{slope}_{T_i}(x) \leq \alpha_i$ for every $x \geq 0$, $i \in [n]$. Let $\alpha = \max_{i \in [n]} \alpha_i$ and $\mu = \max_{i \in [n]} \mu_i$. Finally, let q denote the number of dominating sets over the elements in $[n]$. Then*

$$\mathbf{E}[q] \leq \left(\sum_{i \in [n]} \mu_i \cdot \sum_{i \in [n]} \alpha_i \right) + 1 \leq \alpha \mu n^2 + 1.$$

PROOF. We use an approach similar to the proof of Theorem 2. The bounds that we can prove, however, are slightly better, as we can exploit the Markovian properties of the exponential distribution lower-bounding the long tailed distributions under consideration.

Let $m = 2^n$ and let S_1, \dots, S_m denote the sequence of all subsets of $[n]$ listed in non-decreasing order of their weights. Fix $u \in [m]$,

$u \geq 2$. For every $v \in [u-1]$, define $X_v = S_u \setminus S_v$ and $Y_v = S_v \setminus S_u$. If S_u is a dominating set, then the expected increase in profit is $\mathbf{E}[\Delta_u | \Delta_u > 0]$ corresponding to

$$\begin{aligned} & \mathbf{E} \left[\min_{v \in [u-1]} \left(\sum_{i \in X_v} p_i - \sum_{i \in Y_v} p_i \right) \middle| \min_{v \in [u-1]} \left(\sum_{i \in X_v} p_i - \sum_{i \in Y_v} p_i \right) > 0 \right] \\ &= \int_0^\infty \Pr \left[\forall v: \sum_{i \in X_v} p_i - \sum_{i \in Y_v} p_i \geq t \middle| \forall v: \sum_{i \in X_v} p_i \geq \sum_{i \in Y_v} p_i \right] dt. \end{aligned}$$

Let $k = |S_u|$ be the number of element in S_u . W.l.o.g., assume $S_u = [k]$. Observe that $X_v \subseteq [k]$ and $Y_v \subseteq [n] \setminus [k]$, for every $v \in [u-1]$. Our next goal is to isolate the random variables p_1, \dots, p_k . For this purpose, we partition the set $[u-1]$ into disjoint groups G_1, \dots, G_k satisfying the following property: $\forall j \in [k]: v \in G_j \Rightarrow j \in X_v$. Let $E_j(t)$ denote the event $\forall v \in G_j: \sum_{i \in X_v} p_i \geq \sum_{i \in Y_v} p_i + t$. Then

$$\begin{aligned} \mathbf{E}[\Delta_u | \Delta_u > 0] &= \int_0^\infty \Pr \left[\bigwedge_{i \in [k]} E_i(t) \middle| \bigwedge_{i \in [k]} E_i(0) \right] dt \\ &= \int_0^\infty \prod_{j=1}^k \Pr \left[E_j(t) \middle| \bigwedge_{i=1}^{j-1} E_i(t) \wedge \bigwedge_{i \in [k]} E_i(0) \right] dt. \end{aligned}$$

Now fix some $j \in [k]$ and let us assume that the values of all random variables except for p_j are fixed as well. Define

$$A_j = \max_{v \in G_j} \left(\sum_{i \in Y_v} p_i - \sum_{i \in X_v \setminus \{j\}} p_i \right).$$

Notice that A_j is independent of p_j and, therefore, A_j is also fixed. This way, the expression $E_j(t)$ is equivalent to the expression $p_j \geq A_j + t$. Furthermore, as $E_j(0)$ corresponds to $p_j \geq A_j$, the expression $\bigwedge_{i=1}^{j-1} E_i(t) \wedge \bigwedge_{i \in [k]} E_i(0)$ is equivalent to $p_j \geq A'_j$, for some $A'_j \geq A_j$. Consequently,

$$\begin{aligned} \Pr \left[E_j(t) \middle| \bigwedge_{i=1}^{j-1} E_i(t) \wedge \bigwedge_{i \in [n]} E_i(0) \right] &= \Pr [p_j \geq A_j + t | p_j \geq A'_j] \\ &\geq \Pr [p_j \geq A_j + t | p_j \geq A_j]. \end{aligned}$$

Now recall that T_j is the tail function for the random variable p_j . Hence, $\Pr [p_j \geq A_j + t | p_j \geq A_j] = \frac{T_j(A_j+t)}{T_j(A_j)}$. For the exponential distribution with parameter α , it holds $\frac{T_j(A_j+t)}{T_j(A_j)} = T_j(t) = e^{-\alpha t}$. Here the first equality corresponds to the so-called Markovian or memoryless property of the exponential distribution. For other long-tailed distributions we need a slightly more complicated calculation. Recall, for all $i \in [n]$, $x \in \mathbb{R}_{\geq 0}$, we assume slope $_{T_i}(x) = -[\ln(T_i(x))]' \leq \alpha_i$. This yields

$$\begin{aligned} \ln \left(\frac{T_j(x+t)}{T_j(x)} \right) &= \ln(T_j(x+t)) - \ln(T_j(x)) \\ &\geq (\ln(T_j(x)) - \alpha_j t) - \ln(T_j(x)) = -\alpha_j t \end{aligned}$$

so that $\Pr [p_j \geq A_j + t | p_j \geq A_j] = \frac{T_j(A_j+t)}{T_j(A_j)} \geq e^{-\alpha_j t}$, regardless of the outcome of A_j . Putting all together,

$$\mathbf{E}[\Delta_u | \Delta_u > 0] \geq \int_0^\infty \prod_{j \in [k]} e^{-\alpha_j t} dt \geq \int_0^\infty \prod_{j \in [n]} e^{-\alpha_j t} dt = \frac{1}{\sum_{j \in [n]} \alpha_j},$$

for every $u \in \{2, \dots, n\}$. On the other hand, $\mathbf{E}[P_m] = \sum_{i \in [n]} \mu_i$. Thus, analogous to the proof of Theorem 2, we are now able to

bound the expected number of dominating sets by

$$\mathbf{E}[q] \leq 1 + \frac{\mathbf{E}[P_m]}{\min_{u=2}^m (\mathbf{E}[\Delta_u | \Delta_u > 0])} \leq 1 + \sum_{i \in [n]} \alpha_i \cdot \sum_{j \in [n]} \mu_j.$$

This completes the proof of Theorem 4. \square

3.2 Applications

In this section, we illustrate the power of Theorem 4 by investigating $\mathbf{E}[q]$, the expected number of dominating sets, for some specific long-tailed probability distributions. Recall that the expected running time of the Nemhauser/Ullmann algorithm is $O(nq(n))$, where $q(n)$ is a non-decreasing upper bound on $\mathbf{E}[q]$. In order to simplify the presentation of the results, let us assume that all profits follow the same distribution.

COROLLARY 5. *If profits are chosen according to the exponential distribution, then $\mathbf{E}[q] = O(n^2)$.*

This result is tight. In Section 4 we show a corresponding lower bound. Observe that the result for the exponential distribution does not depend on the parameter of this distribution as the mean μ is reciprocal of the slope α , regardless of the choice for the parameter of this distribution. This is slightly different in the case of the Pareto distribution. For the Pareto distribution with parameter $a > 1$, $\alpha = a$ and $\mu = \frac{1}{a-1}$. This gives the following upper bound on the running time.

COROLLARY 6. *If profits are chosen according to the Pareto distribution with parameter $a > 1$ then $\mathbf{E}[q] = O(\frac{a}{a-1} n^2)$.*

Observe that the Pareto distribution has finite mean only for parameter $a > 1$ so that our proof technique works only for $a > 1$.

4. LOWER BOUND

In this section we prove a lower bound for the number of dominating sets for continuous distributions with non-increasing density functions. This result shows that our upper bound for the exponential distribution is tight.

THEOREM 7. *Suppose profits are drawn independently at random according to a continuous probability distribution with non-increasing density function $f: \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}$. Then there is a vector of weights w_1, \dots, w_n for which $\mathbf{E}[q] = \Omega(n^2)$.*

PROOF. If the density function is non-increasing, then the distribution function $F: \mathbb{R}_{\geq 0} \rightarrow [0, 1]$ is concave as $F' = f$. Furthermore, $F(0) = 0$. Observe that such a function is sub-additive, that is, $F(a+b) \leq F(a) + F(b)$, for every $a, b \geq 0$. This is the crucial property that we need in the following analysis.

Let $w_i = 2^i$ be the weight for the i th item and p_i its profit. For every $j \in [n]$, define $P_j = \sum_{i=1}^j p_i$. Consider the sequence $\mathcal{S}(j-1)$ of dominating sets for items $1, \dots, j-1$ (see Section 1.1). Notice, that all these sets have weight less than 2^j . All sets containing item j , however, have weight at least 2^j and, therefore, these sets cannot dominate the sets in $\mathcal{S}(j-1)$. Hence $\mathcal{S}(j)$ contains all sets from $\mathcal{S}(j-1)$. Furthermore, those sets $S \in \mathcal{S}(j-1)$ with profit $p(S) \in (P_{j-1} - p_j, P_{j-1}]$ create new dominating sets in $\mathcal{S}(j)$ with profit $p(S) + p_j > P_{j-1}$.

For any given $\alpha > 0$, let X_α^j be the number of dominating sets in $\mathcal{S}(j)$ with profit at least $P_j - \alpha$, not counting the last set $[j]$ in this sequence. By induction we show $\mathbf{E}[X_\alpha^j] \geq F(\alpha)j$. Clearly,

$\mathbf{E}[X_\alpha^1] = F(\alpha)$. For $j > 1$, it holds

$$\begin{aligned} \mathbf{E}[X_\alpha^j] &= \Pr[p_j \leq \alpha] \left(\mathbf{E}[X_{p_j}^{j-1}] + \mathbf{E}[X_{\alpha-p_j}^{j-1}] + 1 \right) \\ &\quad + \Pr[p_j > \alpha] \mathbf{E}[X_\alpha^{j-1}] \\ &\geq \Pr[p_j \leq \alpha] (F(p_j)(j-1) + F(\alpha-p_j)(j-1) + 1) \\ &\quad + \Pr[p_j > \alpha] F(\alpha)(j-1) \\ &\stackrel{(*)}{\geq} \Pr[p_j \leq \alpha] (F(\alpha)(j-1) + 1) + \Pr[p_j > \alpha] F(\alpha)(j-1) \\ &= F(\alpha)(j-1) + F(\alpha) = F(\alpha)j, \end{aligned}$$

where inequality (*) follows from $F(a) + F(b) \geq F(a+b)$. Now let $Y_j = |\mathcal{S}(j)| - |\mathcal{S}(j-1)|$ denote the number of new dominating sets in $\mathcal{S}(j)$. Observe that $Y_j = X_{p_j}^{j-1} + 1$. The additive 1 is due to the fact that the set $[j-1]$ is not counted in $X_{p_j}^{j-1}$ but yields a new set in $\mathcal{S}(j)$. Since p_j and X_α^{j-1} are independent, we get $\mathbf{E}[Y_j] = \mathbf{E}[X_{p_j}^{j-1} + 1]$. Furthermore, the number of dominating sets in $\mathcal{S}(n)$ is $q = \sum_{j=1}^n Y_j$ and, therefore,

$$\begin{aligned} \mathbf{E}[q] &= \sum_{j=1}^n \mathbf{E}[Y_j] = \sum_{j=1}^n \mathbf{E}[X_{p_j}^{j-1} + 1] \\ &\geq \sum_{j=1}^n \mathbf{E}[F(p_j)(j-1) + 1] \geq \sum_{j=1}^n \mathbf{E}[F(p_j)]j. \end{aligned}$$

In order to evaluate $\mathbf{E}[F(p_j)]$, we need to examine the distribution function of the random variable $F(p_j)$. Let $F(\cdot)$ denote this distribution function. Recall that p_j is a random variable with distribution function F .

$$F(x) = \Pr[F(p_j) \leq x] = \Pr[p_j \leq F^{-1}(x)] = F(F^{-1}(x)) = x.$$

Thus $F(p_j)$ is uniformly distributed in $[0, 1]$ so that $\mathbf{E}[F(p_j)] = \frac{1}{2}$. Consequently, $\mathbf{E}[q] \geq \frac{1}{2} \sum_{j=1}^n j = \Omega(n^2)$. \square

The theorem shows that our analysis of the expected number of dominating sets for the exponential distribution is tight. The same is true for all long-tailed distributions with finite mean and non-increasing density function. For the uniform distribution, however, lower and upper bound deviate by a factor $\Theta(n)$. Experimental results let us believe that the lower bound is tight and the truth for this distribution is $\Theta(n^2)$ as well.

5. GENERAL DISTRIBUTIONS

In this section, we extend our result towards general, continuous distributions over $\mathbb{R}_{\geq 0}$ with finite mean. The following theorem shows that the expected number of dominating sets increases only linearly with the maximum expected profit and the maximum density over all items. In Section 6 we show how this result can be generalized towards discrete probability distributions.

THEOREM 8. *For every $i \in [n]$, let profit p_i be a non-negative random variable with density function $f_i : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}$. Suppose $\mu \geq \max_{i \in [n]} (\mathbf{E}[p_i])$ and $\phi \geq \max_{i \in [n]} (\max_{x \in \mathbb{R}_{\geq 0}} f_i(x))$. Then the expected number of dominating sets is $\mathbf{E}[q] = O(\phi\mu^4)$.*

PROOF. Unfortunately, the analysis presented in the previous section fails for distributions with short tails. In particular, the idea to lower-bound the increase in profit at every dominating set does not work for short-tailed distributions. In fact, one can define a collection of short-tailed distributions for which the expected increase in profit at some sets, if they become dominating sets, is arbitrarily

small. The point is, however, that those sets are very unlikely to become dominating sets. This property needs to be exploited in the following analysis.

Let $T_i(x) = \int_x^\infty f_i(t)dt$ denote the tail function for the profit of item $i \in [n]$. For each p_i we define an auxiliary random variable $x_i = T_i(p_i)$. Observe that the x_i 's are uniformly distributed over $[0, 1]$. Furthermore, we introduce a cascade of events on which we want to condition. For $k \geq 0$, let X_k denote the event $\forall i \in [n] : x_i \geq 2^{-k}/n$. Observe that $X_{k-1} \subseteq X_k$, for every $k \geq 1$. By conditioning on this cascade of events, we obtain the following upper bound on the expected number of dominating sets:

$$\begin{aligned} \mathbf{E}[q] &= \Pr[X_0] \cdot \mathbf{E}[q | X_0] + \sum_{k=1}^{\infty} \Pr[X_k \wedge \neg X_{k-1}] \cdot \mathbf{E}[q | X_k \wedge \neg X_{k-1}] \\ &\leq \mathbf{E}[q | X_0] + \sum_{k=1}^{\infty} \Pr[\neg X_{k-1}] \cdot \mathbf{E}[q | X_k \wedge \neg X_{k-1}] \\ &\leq \mathbf{E}[q | X_0] + \sum_{k=1}^{\infty} \frac{\mathbf{E}[q | X_k \wedge \neg X_{k-1}]}{2^{k-1}}, \end{aligned} \quad (2)$$

where the last inequality follows because

$$\Pr[\neg X_{k-1}] = \Pr[\exists i \in [n] : x_i < 2^{-k+1}/n] \leq 2^{-k+1}.$$

Unfortunately, placing conditions in such a direct way does not yield longer tails but shorter tails as conditioning on X_k simply cuts the tail function at position $T_i^{-1}(2^{-k}/n)$. The following trick avoids this kind of unwanted effects by *masking* the short tails. For a subset $M \subseteq [n]$, let $q(M)$ denote the number of dominating sets over the items in M . For every $k \geq 0$ and $M \subseteq [n]$, we define an auxiliary random variable $q_k(M)$ as follows.

$$q_k(M) = \begin{cases} q(M) & \text{if } X_k, \\ 0 & \text{otherwise.} \end{cases}$$

For $q_k(M)$ we also write q_k . The variables $q_k(M)$ masks dominating sets in case of $\neg X_k$. The following lemma shows that this masking technique enables us to get rid of conditional probabilities so that our analysis for long tails can be applied to estimate the unconditioned q_k variables subsequently.

LEMMA 9. *For every $k \geq 0$, $\mathbf{E}[q] \leq \sum_{k=0}^{\infty} 2^{-k+5} \mathbf{E}[q_k(M_k)]$, for some $M_k \subseteq [n]$.*

We want to remark that the parameterization of the variables q_k with the sets M_k is only a technicality that can be dropped when assuming that removing an element does not increase the expected number of dominating sets, i.e., $\forall M \subseteq [n] : \mathbf{E}[q([n])] \geq \mathbf{E}[q(M)]$.

PROOF. Here we only sketch the ideas. A full proof of this lemma with all details can be found in the Appendix. First we rewrite Equation 2 in terms of the q_k variables, that is,

$$\mathbf{E}[q] \leq \mathbf{E}[q_0 | X_0] + \sum_{k=1}^{\infty} \frac{\mathbf{E}[q_k | X_k \wedge \neg X_{k-1}]}{2^{k-1}}.$$

Next we give upper bounds for $\mathbf{E}[q_0 | X_0]$ and $\mathbf{E}[q_k | X_k \wedge \neg X_{k-1}]$, and use them in this equation. The first term can be estimated by

$$\mathbf{E}[q_0 | X_0] \leq \frac{\mathbf{E}[q_0]}{\Pr[X_0]} \leq 4\mathbf{E}[q_0].$$

Analogously, the second term can be estimated by

$$\mathbf{E}[q_k | X_k \wedge \neg X_{k-1}] \leq \frac{\mathbf{E}[q_k | \neg X_{k-1}]}{\Pr[X_k | \neg X_{k-1}]} \leq 4\mathbf{E}[q_k | \neg X_{k-1}].$$

Finally, we need to get rid of the conditioning on $\neg X_{k-1}$. This condition states that at least one of the profit variables has a very large value. Roughly speaking, only one variable will be affected by this condition. The idea is now to make use of the fact that every individual profit variable can increase the number of dominating sets only by a factor of two. Using this idea in a formal way, we can show $\mathbf{E}[q_k | \neg X_{k-1}] \leq 4\mathbf{E}[q_k(M_k)]$, for some $M_k \subseteq [n]$, so that $\mathbf{E}[q_k | X_k \wedge \neg X_{k-1}] \leq 16\mathbf{E}[q_k(M)]$. Now substituting these bounds back into the upper bound on $\mathbf{E}[q]$ yields the lemma. \square

Next we apply our analysis for the long tailed distributions to the q_k variables.

LEMMA 10. *For every $k \geq 0$, $\mathbf{E}[q_k] \leq \min\{n^3 2^k \phi \mu + n + 1, 2^n\}$.*

PROOF. The bound $\mathbf{E}[q_k] \leq 2^n$ holds trivially because there are at most 2^n different subsets over n elements. The bound $\mathbf{E}[q_k] \leq n^3 2^k \phi \mu + n + 1$ can be shown with the help of Theorem 4. Define $B_{i,k} = T_i^{-1}(2^{-k}/n)$. Remember that q_k counts only dominating sets under X_k . Therefore, the behavior of the tail function for values larger than $B_{i,k}$ is irrelevant for q_k and we can modify the tail function for values larger than $B_{i,k}$ without affecting q_k . In fact, changing $T_i(x)$, for $x > B_{i,k}$, enables us to bound the slope of the tail function as needed for the application of Theorem 4.

Consider the following variants of our tail function. We cut the tail functions T_i at position $B_{i,k}$ and replace the original, possibly short tails by long, exponential tails. For $i \in [n]$ and $k \geq 0$, define

$$T_{i,k}(t) = \begin{cases} T_i(t) & \text{if } t \leq B_{i,k}, \\ \exp(-\phi n(t - B_{i,k})) / (n 2^k) & \text{if } t > B_{i,k}. \end{cases}$$

The slope of this tail function can be bounded as follows. For $t \leq B_{i,k}$,

$$\text{slope}_{T_{i,k}}(t) = [-\ln(T_{i,k}(t))] = \frac{-[T_i(t)]'}{T_i(t)} \leq \phi n 2^k$$

because $-[T_i(t)]' = f_i(t) \leq \phi$ and $T_i(t) \geq 2^{-k}/n$ since $t \leq B_{i,k} = T_i^{-1}(2^{-k}/n)$. The same upper bound holds also for $t > B_{i,k}$ since, in this case,

$$\text{slope}_{T_{i,k}}(t) = \left[-\ln \left(\frac{\exp(-\phi n(t - B_{i,k}))}{n 2^k} \right) \right]' = \phi n \leq \phi n 2^k.$$

Furthermore, observe that the expected maximum profit of an item under the tail function $T_{i,k}$ is at most $\mu + (\phi n 2^k)^{-1}$ because the added exponential tail increases the original mean value μ by at most $(\phi n)^{-1} (n 2^k)^{-1}$. Consequently, applying Theorem 4 with $\alpha_k = \phi n 2^k$ and $\mu_k = \mu + (\phi n 2^k)^{-1}$ yields $\mathbf{E}[q_k] \leq \alpha_k \mu_k n^2 + 1 = \phi \mu n^3 2^k + n + 1$. \square

Now we can complete our calculation for the upper bound on $\mathbf{E}[q]$. Combining Lemma 9 and 10 yields

$$\begin{aligned} \mathbf{E}[q] &\leq \sum_{k=0}^{\infty} 2^{-k+5} \min\{\phi \mu n^3 2^k + n + 1, 2^n\} \\ &\leq 32 \sum_{k=0}^n (\phi \mu n^3 + 2^{-k}(n+1)) + 32 \sum_{k=n+1}^{\infty} 2^{n-k} \\ &\leq 32(\phi \mu n^4 + 2n + 3). \end{aligned}$$

Finally observe, $\phi \mu \geq \frac{1}{2}$, for every non-negative continuous distribution. Thus $\mathbf{E}[q] = O(\phi \mu n^4)$, which completes the proof of Theorem 8.

6. DISCRETE DISTRIBUTIONS

In this section we generalize our results towards discrete probability distributions, that is, we assume that profits are randomly drawn non-negative integers. In this case, we can prove a trade-off ranging from polynomial to pseudo-polynomial running time, depending on the degree of randomness of the specified instances. The proofs rely on the same methods that we used for continuous distributions. So we will point out only those parts of the proofs that need to be changed. First, we consider long-tailed and then general discrete distributions.

6.1 Long-tailed discrete distributions

Let us assume that profits are chosen independently at random according to possibly different long-tailed distributions with finite mean. For $i \in [n]$, let profit p_i be a random variable with tail function $T_i : \mathbb{N}_0 \rightarrow [0, 1]$, that is, for every $t \in \mathbb{N}_0$, $\Pr[p_i \geq t] = T_i(t)$.

THEOREM 11. *Let α be an appropriate positive term satisfying $T_i(t+1)/T_i(t) \geq e^{-\alpha}$ for all $i \in [n]$ and $t \in \mathbb{N}_0$. Suppose $\mu \geq \max_{i \in [n]} (\mathbf{E}[p_i])$. Let q denote the number of dominating sets over the elements in $[n]$. Then $\mathbf{E}[q] \leq \mu n(1 - e^{-\alpha n}) + 1 \leq \mu \alpha n^2 + 1$.*

For the application of this theorem, it makes sense to assume that the probability distributions “scale” with the number of items. For example, consider the following discrete variant of the exponential distribution. Let $T_i(t) = e^{-\alpha t}$ ($t \geq 0$, $i \in [n]$) with $\alpha = \alpha(n)$ being a function in n , e.g., $\alpha(n) = \frac{1}{n}$. The mean μ of this distribution grows like $\frac{1}{\alpha(n)}$. Thus $\mathbf{E}[q] = O(n^2)$ under the discrete exponential distribution, regardless of the choice of α . In other words, we obtain the same bound as in the continuous case.

PROOF. We can adapt the proof of Theorem 4 towards discrete long tailed distributions.

$$\begin{aligned} \mathbf{E}[\Delta_u | \Delta_u \geq 1] &= \mathbf{E} \left[\min_{v \in [u-1]} \left(\sum_{i \in X_v} p_i - \sum_{i \in Y_v} p_i \right) \middle| \min_{v \in [u-1]} \left(\sum_{i \in X_v} p_i - \sum_{i \in Y_v} p_i \right) \geq 1 \right] \\ &= \sum_{t=1}^{\infty} \Pr \left[\forall v : \sum_{i \in X_v} p_i \geq \sum_{i \in Y_v} p_i + t \mid \forall v : \sum_{i \in X_v} p_i \geq \sum_{i \in Y_v} p_i + 1 \right]. \end{aligned}$$

We define the groups G_1, \dots, G_k as in the proof of Theorem 4. Let $E_j(t)$ denote the event $\forall v \in G_j : \sum_{i \in X_v} p_i \geq \sum_{i \in Y_v} p_i + t$. Then

$$\begin{aligned} \mathbf{E}[\Delta_u | \Delta_u \geq 1] &= \sum_{t=1}^{\infty} \Pr \left[\bigwedge_{i \in [k]} E_i(t) \mid \bigwedge_{i \in [k]} E_i(1) \right] \\ &= \sum_{t=1}^{\infty} \prod_{j=1}^k \Pr \left[E_j(t) \mid \bigwedge_{i=1}^{j-1} E_i(t) \wedge \bigwedge_{i \in [k]} E_i(1) \right] \\ &\geq \sum_{t=1}^{\infty} \prod_{j \in [k]} e^{-\alpha(t-1)} \geq \sum_{t=1}^{\infty} e^{-\alpha n(t-1)} = \frac{1}{1 - e^{-\alpha n}}. \end{aligned}$$

We used the fact that $T_j(A_j + t)/T_j(A_j + 1) \geq e^{-\alpha(t-1)}$ for all $j \in [k]$, $t \geq 1$ where A_j is defined as in the proof of Theorem 4. Notice, that the expected increase in profit at dominating sets is lower bounded by 1 which reflects the integrality of profits. Now we can bound the number of dominating sets by

$$\begin{aligned} \mathbf{E}[q] &= 1 + \sum_{u=2}^m \Pr[\Delta_u > 0] \leq 1 + \frac{\mathbf{E}[P_m]}{\min_{u=2}^m (\mathbf{E}[\Delta_u | \Delta_u > 0])} \\ &\leq \mu n(1 - e^{-\alpha n}) + 1. \end{aligned}$$

Applying the inequality $1 - e^{-x} \leq x$ yields $\mathbf{E}[q] \leq \mu \alpha n^2 + 1$. \square

6.2 General discrete distributions

In this section, we analyze the number of dominating sets when profits are chosen according to general discrete probability distributions over \mathbb{N} with finite mean. For every $i \in [n]$, we assume that p_i is a positive random variable with probability function $f_i : \mathbb{N} \rightarrow [0, 1]$, i.e., $f_i(t) = \Pr[p_i = t]$.

THEOREM 12. *Suppose $\pi = \max_{i \in [n]} (\max_{x \in \mathbb{N}} (f_i(x)))$ and $\mu \geq \max_{i \in [n]} (\mathbf{E}[p_i])$. Then the expected number of dominating sets is $\mathbf{E}[q] = O(\mu n^2 (1 - e^{-\pi n^2})) = O(\mu \pi n^4)$.*

In fact, the term $1 - e^{-\pi n^2}$ in the upper bound translates into a pseudo-polynomial bound if the randomness in the specification goes to zero. For example, assume that for each item an adversary specifies an interval from which the profit of this item is drawn uniformly at random. Let M denote the maximum profit that can be drawn for any item and ℓ the minimum interval length over all items. Set $\mu = M$ and $\pi = \frac{1}{\ell}$ so that $\mathbf{E}[q] = O(Mn^2(1 - e^{-n^2/\ell}))$. Now, if $\ell = \Theta(M)$ then this upper bound simplifies to $O(n^4)$ because $1 - e^{-x} \leq x$, for all $x \in \mathbb{R}$. However, if $\ell = \Theta(1)$, then we are left with the pseudo-polynomial upper bound $\mathbf{E}[q] = O(Mn^2)$.

PROOF. We adapt the proof of Theorem 8. Let $T_i : \mathbb{N} \rightarrow [0, 1]$ be the tail function of p_i . To simplify the analysis, we introduce auxiliary random variables x_1, \dots, x_n . These variables are drawn independently and uniformly over the interval $[0, 1]$. Now we assume that the p_i 's are generated from the x_i 's by setting $p_i = \max\{j \in \mathbb{N} : T_i(j) \geq x_i\}$. In this way, $\Pr[p_i \geq t] = T_i(t)$, so that this is a proper way to generate the profits in order to obtain the same distribution as described in the theorem. Define auxiliary random variables $q_k(M)$ as before. Lemma 9 from the proof of Theorem 8 transfers directly to discrete probability distributions. For the convenience of the reader we state it here again.

LEMMA 13. *For every $k \geq 0$, $\mathbf{E}[q] \leq \sum_{k=0}^{\infty} 2^{-k+5} \mathbf{E}[q_k(M_k)]$, for some $M_k \subseteq [n]$.*

Next we prove the discrete counterpart of Lemma 10.

LEMMA 14. *For $k \geq 0$, $\mathbf{E}[q_k] \leq \min\{\mu n(1 - e^{-\pi 2^k n^2}) + n + 1, 2^n\}$.*

PROOF. The bound $\mathbf{E}[q_k] \leq 2^n$ holds trivially because there are at most 2^n different subsets over n elements. The other bound can be shown with the help of Theorem 11. In order to apply this Theorem, we need to bound the ratio of $T_{i,k}(t+1)$ and $T_{i,k}(t)$ for all $t \in \mathbb{N}$. Define $B_{i,k} = \min\{j \in \mathbb{N} : T_i(j) < 2^{-k}/n\}$. Consider the following variants of our tail function. We cut the tail functions T_i at position $B_{i,k}$ and replace the original, possibly short tails by long, exponential tails. For $i \in [n]$ and $k \geq 0$, define

$$T_{i,k}(t) = \begin{cases} T_i(t) & \text{if } t < B_{i,k}, \\ \exp(-\pi n(t - B_{i,k}))/n2^k & \text{if } t \geq B_{i,k}. \end{cases}$$

Figure 3 illustrates the situation.

For all $t \leq B_{i,k} - 1$, $T_{i,k}(t+1) \geq 2^{-k}/n$ and $T_{i,k}(t) - T_{i,k}(t+1) \leq \pi$. In particular, $T_{i,k}(B_{i,k} - 1) - T_{i,k}(B_{i,k}) \leq \pi$, because $T_{i,k}(B_{i,k}) \geq T(B_{i,k})$. Therefore $T_{i,k}(t) \leq T_{i,k}(t+1) + 2^k n \cdot T_{i,k}(t+1) \cdot \pi$ and

$$\frac{T_{i,k}(t+1)}{T_{i,k}(t)} \geq \frac{1}{1 + 2^k n \pi} \geq e^{-2^k n \pi}.$$

The same lower bound holds also for $t \geq B_{i,k}$ since, in this case,

$$\frac{T_{i,k}(t+1)}{T_{i,k}(t)} = e^{-n\pi} \geq e^{-2^k n \pi}.$$

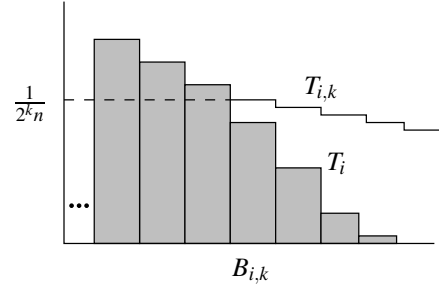


Figure 3: We cut the tail T_i at position $B_{i,k}$ and replace this part by a long tail $T_{i,k}$.

Furthermore, observe that the expected maximum profit of an item under the tail function $T_{i,k}$ is at most $\mu + 1/(2^k n(1 - e^{-\pi n}))$, because the added exponential tail increases the original mean value μ by at most

$$\max_i \sum_{t=B_{i,k}}^{\infty} \frac{1}{2^k n} e^{-\pi n(t-B_{i,k})} = \frac{1}{2^k n} \cdot \frac{1}{(1 - e^{-\pi n})}.$$

Applying Theorem 11 with parameters $\mu_k = \mu + 1/(2^k n(1 - e^{-\pi n}))$ and $\alpha_k = \pi n 2^k$ yields

$$\begin{aligned} \mathbf{E}[q_k] &\leq \mu_k n(1 - e^{-\alpha_k n}) + 1 \\ &= \mu n(1 - e^{-\pi n 2^k}) + \frac{1 - e^{-\pi n 2^k}}{2^k(1 - e^{-\pi n})} + 1 \\ &\leq \mu n(1 - e^{-\pi n 2^k}) + n + 1. \end{aligned}$$

Thus, Lemma 14 is shown. \square

Now we can complete our calculation for the upper bound on $\mathbf{E}[q]$. Combining Lemma 14 and 13 yields

$$\begin{aligned} \mathbf{E}[q] &\leq \sum_{k=0}^{\infty} 2^{-k+5} \min\left\{\mu n(1 - e^{-\pi n 2^k}) + n + 1, 2^n\right\} \\ &\leq 32 \left(\sum_{k=0}^n \mu n \frac{(1 - e^{-\pi n 2^k})}{2^k} + \frac{n+1}{2^k} \right) + 32 \sum_{k=n+1}^{\infty} 2^{n-k} \\ &\leq 32(\mu n^2(1 - e^{-\pi n^2}) + 2n + 3). \end{aligned}$$

Finally, we need to show that $2n + 3 = O(\mu n^2(1 - e^{-\pi n^2}))$. As the considered distributions are positive, $\mu \geq 1$ and $\pi \mu \geq \frac{1}{2}$ gives $(1 - e^{-\pi n^2}) \geq (1 - e^{-n^2/2\mu})$. The inequality $x + e^{-n^3x/2} \leq 1$ holds for all $0 \leq x \leq 0.5$ and $n \geq 2$. Substituting $x := 1/(\mu n)$ yields $1 \leq \mu n(1 - e^{-n^2/2\mu})$. Thus Theorem 12 is shown.

7. ACKNOWLEDGEMENTS

We like to thank Micah Adler, Kurt Mehlhorn and Uli Meyer for inspiring discussions, George Lueker for providing useful information about previous work as well as Alan Frieze for helpful comments on an earlier version of this paper.

8. REFERENCES

- [1] E. Balas and E. Zemel. An algorithm for large zero-one knapsack problems. *Operations Research*, Vol. 28, pp. 1130-1154, 1980.
- [2] K. H. Borgwardt and J. Brzank. Average Saving Effects in Enumerative Methods for Solving Knapsack Problems. *J. of Complexity* 10, pp. 129-141, 1994.

- [3] M. Coster, A. Joux, B. LaMacchia, A. Odlyzko, C. P. Schnorr, J. Stern. Improved Low-Density Subset Sum Algorithms. *J. of Computational Complexity*, 111–128, 1992.
- [4] G. B. Dantzig. Discrete Variable Extremum Problems. *Operations Research*, Vol 5, pp. 266-277, 1957.
- [5] G. D'Atri and C. Puech. Probabilistic Analysis of the Subset-Sum Problem. *Discrete Applied Math.*, Vol 4, 329–334, 1982.
- [6] M. E. Dyer and A. M. Frieze. Probabilistic Analysis of the Multidimensional Knapsack Problem. *Mathematics of Operations Research* 14(1), 162-176, 1989.
- [7] D. Fayard and G. Plateau. An algorithm for the solution of the 0-1 knapsack problem. *Computing*, Vol. 28, pp. 269-287, 1982.
- [8] A. M. Frieze. On the Lagarias-Odlyzko algorithm for the Subset-Sum Problem. *SIAM J. Comput.* 15(2), 536-539, 1986.
- [9] M. R. Garey and D. S Johnson. *Computers and Intractability: A Guide to the Theory of NP-completeness*. Freeman, 1979.
- [10] S. Martello, D. Pisinger and P. Toth. New Trends in Exact Algorithms for the 0/1 Knapsack Problem. *European Journal of Operational Research*, Vol. 123, pp. 325-332, 2000.
- [11] S. Martello and P. Toth. *Knapsack Problems – Algorithms and Computer Implementations*. Wiley, 1990.
- [12] A. Goldberg and A. Marchetti-Spaccamela. On Finding the Exact Solution to a Zero-One Knapsack Problem. *Proceedings of the 16th ACM Symposium on Theory of Computing (STOC)*, 359–368, 1984.
- [13] R. Impagliazzo and M. Naor. Efficient Cryptographic Schemes Provably as Secure as Subset Sum. *J. of Cryptology* 9(4), 199–216, 1996.
- [14] G. S. Lueker. On the Average Difference Between the Solutions to Linear and Integer Knapsack Problems. *Applied Probability - Computer Science, the Interface*, 1, Birkhauser, 1982.
- [15] G. S. Lueker. Average-Case Analysis of Off-Line and On-Line Knapsack Problems. *J. of Algorithms*, 19, 277-305, 1998.
- [16] G. S. Lueker. Exponentially Small Bounds on the Expected Optimum of the Partition and Subset Sum Problems. *Random Structures and Algorithms*, 12, 51-62, 1998.
- [17] G. Nemhauser and Z. Ullmann. Discrete Dynamic Programming and Capital Allocation. *Management Science*, 15(9), 494–505, 1969.
- [18] E. Horowitz and S. Sahni. Computing Partitions with Applications to the Knapsack Problem, *J. of the ACM*, Vol. 21, 277–292, 1974.
- [19] D. Pisinger. *Algorithms for Knapsack Problems*. Ph.D. thesis, DIKU, University of Copenhagen, 1995.
- [20] D. Pisinger and P. Toth. *Knapsack Problems*. In D-Z. Du, P. Pardalos (ed.), *Handbook of Combinatorial Optimization*, vol. 1, Kluwer Academic Publishers, 299-428, 1998.
- [21] D. A. Spielman and Shang-Hua Teng. Smoothed Analysis of Algorithms: Why The Simplex Algorithm Usually Takes Polynomial Time. *Proceedings of the 33rd ACM Symposium on Theory of Computing (STOC)*, 296–305, 2001.
- [22] H.M. Weingartner and D.N. Ness. Methods for the Solution of the Multi-Dimensional 0/1 Knapsack problem. *Operations Research*, Vol. 15, No. 1, 83–103, 1967.

APPENDIX

A. FORMAL PROOF OF LEMMA 9

We have to show $\mathbf{E}[q] \leq \sum_{k=0}^{\infty} 2^{-k+5} \mathbf{E}[q_k]$, for every $k \geq 0$. First we rewrite Equation 2 in terms of the q_k variables, that is,

$$\mathbf{E}[q] \leq \mathbf{E}[q_0 | X_0] + \sum_{k=1}^{\infty} \frac{\mathbf{E}[q_k | X_k \wedge \neg X_{k-1}]}{2^{k-1}}. \quad (3)$$

Next we prove upper bounds for $\mathbf{E}[q | X_0]$ and $\mathbf{E}[q | X_k \wedge \neg X_{k-1}]$. The first term can be estimated as follows.

$$\mathbf{E}[q_0 | X_0] \leq \frac{\mathbf{E}[q_0]}{\Pr[X_0]} \leq 4\mathbf{E}[q_0], \quad (4)$$

where the last equation follows as $\Pr[X_0] = \Pr[\forall i \in [n] : x_i \geq \frac{1}{n}] = (1 - \frac{1}{n})^n \geq \frac{1}{4}$, for $n \geq 2$. In order to bound the second term, we

partition the event $\neg X_{k-1}$ into n disjoint events Z_1, \dots, Z_n with

$$Z_j = \left[x_1, x_2, \dots, x_{j-1} \geq \frac{1}{2^{k-1}n} \wedge x_j < \frac{1}{2^{k-1}n} \right].$$

Then, for $k \geq 1$,

$$\begin{aligned} \Pr[X_k | \neg X_{k-1}] &= \sum_{j=1}^n \Pr[X_k | Z_j] \cdot \Pr[Z_j | \bigcup_{i \in [n]} Z_i] \\ &= \sum_{j=1}^n \frac{1}{2} \left(1 - \frac{1}{2^{k-1}n}\right)^{n-j} \Pr[Z_j | \bigcup_{i \in [n]} Z_i] \\ &\geq \frac{1}{2} \left(1 - \frac{1}{2^{k-1}n}\right)^{n-1} \sum_{j=1}^n \Pr[Z_j | \bigcup_{i \in [n]} Z_i] \\ &= \frac{1}{2} \left(1 - \frac{1}{2^{k-1}n}\right)^{n-1} \geq \frac{1}{4}. \end{aligned}$$

Thus, we get

$$\mathbf{E}[q_k | X_k \wedge \neg X_{k-1}] \leq \frac{\mathbf{E}[q_k | \neg X_{k-1}]}{\Pr[X_k | \neg X_{k-1}]} \leq 4\mathbf{E}[q_k | \neg X_{k-1}]. \quad (5)$$

Finally, we need to get rid of the conditioning on $\neg X_{k-1}$. For all $i \in [n]$, let Y_i denote the event $x_i < 2^{-(k-1)}/n$. As $\neg X_{k-1} = \bigcup_{i \in [n]} Y_i$,

$$\begin{aligned} \mathbf{E}[q_k | \neg X_{k-1}] &= \mathbf{E}\left[q_k \mid \bigcup_{i \in [n]} Y_i\right] \\ &\leq \sum_{j \in [n]} \mathbf{E}[q_k | Y_j] \cdot \Pr\left[Y_j \mid \bigcup_{i \in [n]} Y_i\right]. \end{aligned}$$

Observe that the last estimation would hold with equality if the events Y_1, \dots, Y_n were disjoint. As these events overlap, however, the right hand term slightly overestimates the left hand term. Now $\Pr[Y_j] = 2^{-k+1}/n$ and $\Pr\left[\bigcup_{i \in [n]} Y_i\right] = 1 - (1 - 2^{-k+1}/n)^n \geq 2^{-k}$, so that

$$\Pr\left[Y_j \mid \bigcup_{i \in [n]} Y_i\right] = \frac{\Pr[Y_j]}{\Pr\left[\bigcup_{i \in [n]} Y_i\right]} \leq \frac{2^{-k+1}/n}{2^{-k}} \leq \frac{2}{n}.$$

As a consequence,

$$\begin{aligned} \sum_{j \in [n]} \mathbf{E}[q_k | Y_j] \cdot \Pr\left[Y_j \mid \bigcup_{i \in [n]} Y_i\right] &\leq \max_{j \in [n]} (\mathbf{E}[q_k | Y_j]) \sum_{j \in [n]} \Pr\left[Y_j \mid \bigcup_{i \in [n]} Y_i\right] \\ &\leq 2 \max_{j \in [n]} (\mathbf{E}[q_k | Y_j]). \end{aligned}$$

Let j denote the element maximizing the last expression. Observe that $\mathbf{E}[q_k | Y_j] \leq 2\mathbf{E}[q_k([n] \setminus \{j\})]$ because adding a single item can increase the number of dominating sets at most by a factor of two. (To see this, assume that the item is inserted as last element using the Nemhauser/Ullmann algorithm.) Hence,

$$\mathbf{E}[q_k | \neg X_{k-1}] \leq 2(\mathbf{E}[q_k | Y_j]) \leq 4\mathbf{E}[q_k([n] \setminus \{j\})]. \quad (6)$$

Now substituting the Equations (4), (5), and (6) into Equation (3) yields Lemma 9.