

# On Sums of Independent Random Variables with Unbounded Variance, and Estimating the Average Degree in a Graph

Uriel Feige \*

September 9, 2005

## Abstract

We prove the following inequality: for every positive integer  $n$  and every collection  $X_1, \dots, X_n$  of nonnegative independent random variables that each has expectation 1, the probability that their sum remains below  $n + 1$  is at least  $\alpha > 0$ . Our proof produces a value of  $\alpha = 1/13 \simeq 0.077$ , but we conjecture that the inequality also holds with  $\alpha = 1/e \simeq 0.368$ .

As an example for the use of the new inequality, we consider the problem of estimating the average degree of a graph by querying the degrees of some of its vertices. We show the following threshold behavior: approximation factors above 2 require far less queries than approximation factors below 2. The new inequality is used in order to get tight (up to multiplicative constant factors) relations between the number of queries and the quality of the approximation. We show how the degree approximation algorithm can be used in order to quickly find those edges in a network that belong to many shortest paths.

## 1 A new inequality

For a random variable  $X$ , its typical value may be very different from its mean. In particular, the probability that  $X$  exceeds its mean may be arbitrarily close to 1. In some special cases (e.g., when  $X$  is symmetric around its mean), the probability that  $X$  exceeds its mean is at most  $1/2$ . The purpose of this manuscript is to investigate the probability that  $X$  exceeds its mean when  $X$  is the sum of  $n$  independent random variables. We show that for nonnegative random variables, this probability is bounded away from 1, provided that we give ourselves a little slackness in exceeding the mean.

---

\*Department of Computer Science and Applied Mathematics, the Weizmann Institute, Rehovot 76100, Israel, and Microsoft Research, One Microsoft Way, Redmond, WA 98052-6399, USA. [uriel.feige@weizmann.ac.il](mailto:uriel.feige@weizmann.ac.il), [urifeige@microsoft.com](mailto:urifeige@microsoft.com).

**Theorem 1** *Let  $X_1, \dots, X_n$  be arbitrary nonnegative independent random variables, with expectations  $\mu_1, \dots, \mu_n$  respectively, where  $\mu_i \leq 1$  for every  $i$ . Let  $X = \sum_{i=1}^n X_i$ , and let  $\mu$  denote the expectation of  $X$  (hence,  $\mu = \sum_{i=1}^n \mu_i$ ). Then for every  $\delta > 0$ ,*

$$\Pr[X < \mu + \delta] \geq \min[\delta/(1 + \delta), 1/13] \quad (1)$$

The term  $\delta/(1 + \delta)$  in Theorem 1 is best possible, as one can take  $X_1 = 1 + \delta$  with probability  $1/(1 + \delta)$  and 0 otherwise, and all of the other  $X_i$  as the constant 1. This gives  $\mu_i = 1$  for every  $i$ . For this case  $\Pr[X < \mu + \delta] = \Pr[X_1 = 0] = \delta/(1 + \delta)$ . For large  $\delta$  (e.g.,  $\delta = 1$ ), it is not true that  $\Pr[X \leq \mu + \delta] \geq \delta/(1 + \delta)$ . One can take for every  $i$ ,  $X_i = n + \delta$  with probability  $1/(n + \delta)$  and 0 otherwise. This gives  $\mu_i = 1$  for every  $i$ , implying  $\mu = n$ . For this case  $\Pr[X < n + \delta] = (1 - 1/(n + \delta))^n$ , which is roughly  $1/e$  for large  $n$ .

Based on the two examples above we make the following conjecture.

**Conjecture 1** *In the setting of Theorem 1, for every value of  $\delta$  and  $n$ , one of the two examples above is the worst case for  $\Pr[X < \mu + \delta]$ .*

Conjecture 1, if true, would allow us to replace the constant  $1/13$  by  $1/e$  in Theorem 1. A conjecture very similar in nature to Conjecture 1 is the following conjecture that was suggested by Samuels in [6].

**Conjecture 2** *Let  $X_1, \dots, X_n$  be arbitrary nonnegative independent random variables, with expectations  $\mu_1, \dots, \mu_n$  respectively, where  $\mu_1 \geq \dots \geq \mu_n$ . Then for every  $\lambda > \sum_{j=1}^n \mu_j$  there is some  $1 \leq i \leq n$  such that  $\Pr[\sum_{j=1}^n X_j < \lambda]$  is minimized when the random variables  $X_j$  are distributed as follows:*

- For  $j > i$ ,  $X_j = \mu_j$  with probability 1.
- For  $j \leq i$ ,  $X_j = \lambda - \sum_{k=i+1}^n \mu_k$  with probability  $\frac{\mu_j}{\lambda - \sum_{k=i+1}^n \mu_k}$ , and  $X_j = 0$  otherwise.

The difference between the settings of the two conjectures is that in Conjecture 2 all means are given, whereas in Conjecture 1 only an upper bound on the means is given. The difference in the conclusions of the conjectures is that in Conjecture 2 we may have an arbitrary  $i \in \{1, \dots, n\}$ , whereas Conjecture 1 effectively states that  $i$  is restricted to one of two values,  $i \in \{1, n\}$ .

Samuels (see [6, 7] and references therein) proves Conjecture 2 when  $n$  (the number of random variables) is at most four. The case  $n = 2$  was proved earlier in [2]. When  $n \geq 5$ , Samuels shows that Conjecture 2 is true when  $\lambda \geq (n - 1) \sum_{j=1}^n \mu_j$ . In contrast, in the current paper we are interested in the case (that when put in the framework of Conjecture 2 corresponds to)  $\lambda = \delta + \sum_{j=1}^n \mu_j$  with  $\delta$  fairly small (e.g.,  $\delta = \mu_1$ ).

It may be instructive to consider how some standard probabilistic tools relate to Theorem 1. Consider the case that the  $X_i$  are identically distributed. Then the central limit theorem implies that when  $n$  is large enough,  $X$  approaches

the normal distribution and hence  $Pr[X < \mu]$  approaches  $1/2$ . However, in our Theorem 1 the variables  $X_i$  may depend on  $n$ , and hence  $n$  cannot be thought of as being “large enough” with respect to the  $X_i$  (even if they are i.i.d.). This relates to the fact that we place no bounds on the variance of the  $X_i$ , and hence standard bounds on deviations of random variables from their expectation (such as Chebyshev’s bound, or Chernoff’s bound) are not applicable. The only restriction on the random variables (other than being independent) is their nonnegativity. In particular, this means that  $X$  is nonnegative, and that Markov’s inequality can be used to show that  $Pr[X \leq \mu + \delta] \geq \delta/(\mu + \delta)$ . For the sum of independent identically distributed random variables, this bound tends to 0 as  $n$  grows (unlike the bound in Theorem 1).

In addition to the work by Samuels mentioned above, the author is aware of some other work of nature similar to Theorem 1. There are results surveyed and developed by Siegel [8] that show that under certain conditions the median of the sum of random variables does not exceed the mean. This holds for example for the sum of Bernoulli random variables (if the mean is an integer). The book “How to gamble if you must” by Dubins and Savage [4] analyses strategies for gambling when the goal is to maximize the probability of ending up with a net profit of  $\delta$ . There the strategies are adaptive (next gamble may depend on outcomes of previous gambles) and the gambler may quit once a net profit of  $\delta$  is achieved. One of the main findings of [4] is a set of sufficient conditions under which the strategy of “playing boldly” is optimal. Informally, this strategy tries to reach a net profit of  $\delta$  (taking into account also previous losses) in one gamble. A simple example is the repeated doubling approach to gain one dollar when there are 50/50 odds, in which the gambler first gambles one dollar, and then doubles the gamble until the first win (or until he/she runs out of money). The scenario in Theorem 1 can be viewed as a version of “how to gamble in parallel”, in which  $n$  unbiased gambles with independent outcomes can be placed in parallel in an attempt to reach a net profit of  $\delta$  units, where each gamble is allowed to risk at most one unit. Our results show that when  $\delta$  is small (specifically,  $\delta \leq 1/12$ ), having  $n > 1$  does not lead to higher probability of achieving a net profit of  $\delta$  compared to the case that  $n = 1$ . For larger values of  $\delta$ , there is an advantage of having  $n > 1$ , but regardless of the value of  $n$ , the probability of achieving a net profit of  $\delta$  is bounded away from 1. Regardless of the value of  $\delta$ , it appears to us (though our proof does not actually show this when  $\delta > 1/12$ ) that similar to the “play boldly” principle, the optimal strategy is based on hoping for one successful gamble. (Namely, when  $\delta$  is small only one gamble is nonzero, and when  $\delta$  is large, all gambles are identical, and it suffices for one successful gamble to both reach a profit of  $\delta$  and cover for all the losses in the other gambles.) Despite similarities in the nature of the results, the proof techniques from [4] and [8] do not appear to be applicable to the setting of Theorem 1.

Theorem 1 can in principle be used whenever one is interested in bounding the probability that the sum of independent random variables significantly exceeds its expectation. However, in many contexts the random variables are known to have some additional properties (e.g., bounded variance), and useful

results can also be derived by other means. The application that motivated the development of the inequality (1) is described in Section 2.

## 2 Estimating the average degree

Let  $G(V, E)$  be a graph with  $n$  vertices. A *degree query* specifies a vertex  $v \in V$ , and gets in reply  $d_v$ , the degree of  $v$  in  $G$ . We are interested in estimating  $m = |E|$  by making only degree queries. Equivalently, we wish to estimate the average degree  $d = 2m/n$ . We say that an algorithm provides a  $\rho$  estimation if its output  $d^*$  satisfies

$$d^* \leq d \leq \rho d^*.$$

Naturally, we limit our interest to  $\rho \geq 1$ . As our sampling based algorithms are randomized, there is some probability that their output fails to be a  $\rho$  estimation. We require this failure probability to be at most  $1/3$ . We note that the failure probability can be reduced to an arbitrarily small value  $\delta$ , by repeating the estimation algorithm  $O(\log 1/\delta)$  times independently, and outputting the median of all estimates. Our goal is for given  $\rho$ , to design  $\rho$  estimation algorithms with as few queries as possible, and with failure probability at most  $1/3$ .

Let us note here an observation that helps us to drastically reduce the number of queries in our algorithms. Consider first the case where rather than having an actual graph as input, the input is simply a sequence of integers  $d_1, \dots, d_n$ , with the only restriction that for every  $i$ ,  $0 \leq d_i \leq n$ . (For simplicity of the presentation we allow here values to range up to  $n$ , even though degrees can range only up to  $n - 1$ .) Let  $d = \frac{1}{n} \sum_{i=1}^n d_i$ . We wish to estimate  $d$ . It is not hard to see that for any value  $d_0$  (which one may think of as a large constant independent of  $n$ ),  $\Omega(n/d_0)$  queries are required in order to distinguish between the cases  $d = 0$  and  $d \geq d_0$ . The reason is that it may happen that there are  $d_0$  numbers with value  $n$ , and all other numbers have value 0. If we perform less than  $n/2d_0$  queries, most likely we always get the 0 answer, which is exactly the answers that we would get if  $d = 0$ .

To get estimation algorithms with fewer queries, we shall use the fact that not every sequence  $d_1, \dots, d_n$  is a degree sequence of graphs. For example, if  $d_1 = n - 1$ , then necessarily  $d_i \geq 1$  for all  $i$ . Still, the bad example given above can be modified to show that  $\Omega(n/d_0)$  queries are required in order to distinguish between the cases  $d \leq d_0$  and  $d \geq 2d_0 - O((d_0)^2/n)$ . In the first of these two cases we can have all  $d_i = d_0$ . In the second of these two cases we can have  $d_i = d_0$  for all vertices except for  $d_0$  vertices of degree  $n - 1$ . Hence if we wish to have estimation algorithms with a sublinear (in  $n$ ) number of queries, we need to restrict ourselves to  $\rho \geq 2$ .

There is one more restriction that we introduce. Observe that if  $G$  contains only one edge, one needs  $\Omega(n)$  queries to distinguish this case from  $d = 0$ . To avoid the problem of handling such very sparse graphs (which are often not interesting anyway), we shall assume that  $d \geq d_0$ , for some  $d_0$  that will be a parameter of our estimation algorithms. The reader may think of  $d_0$  as typically

having value at least 1. Hence the estimation algorithm is allowed to output  $d^* = 0$  as an estimation of  $d$  for very sparse graphs, even though the ratio between  $d$  and  $d^*$  is in this case infinite. (The assumption that  $d \geq d_0$  can be replaced by allowing the estimation algorithm to have an additive error of  $d_0$  in addition to the multiplicative error of  $\rho$ .)

As noted above, for  $\rho < 2$  and  $d_0 = 1$ , the number of queries needed by an estimation algorithm might be  $\Omega(n)$ . Our main observation is that for  $\rho > 2$  and for  $d_0 = 1$ , the number of queries in the estimation algorithm drops dramatically, from  $\Omega(n)$  to  $O(\sqrt{n})$ . This result is stated in more technical terms in the following theorem.

**Theorem 2** *For any  $d_0 > 0$ ,  $\epsilon > 0$  and  $\rho = 2 + \epsilon$ , there is a  $\rho$  estimation algorithm for the average degree of a graph that uses  $O(\frac{1}{\epsilon}\sqrt{n/d_0})$  queries, and is applicable to all graphs of average degree at least  $d_0$ .*

In terms of the application of estimating the average degree in the graph, the more interesting part of our upper bound on the number of queries is the term  $\sqrt{n/d_0}$ . The dependency on  $\epsilon$  may be less interesting, especially if one is satisfied with large values of  $\epsilon$ , such as  $\epsilon = 1$ . However, achieving a linear dependency in  $1/\epsilon$  (rather than some polynomial dependency) is the part that uses Theorem 1.

In Section 4 we prove Theorem 2. In Section 5 we show how Theorem 2 can be used in order to obtain Theorem 3, addressing a problem that is studied in [3].

**Theorem 3** *There is a randomized algorithm that runs in time  $O(\frac{mn \log n}{\epsilon\sqrt{c}})$  on graphs with  $n$  vertices and  $m$  edges, and outputs a list of edges that with high probability satisfies:*

1. *Every edge that is on at least  $c$  shortest paths is on the list.*
2. *No edge that is on less than  $(1/2 - \epsilon)c$  shortest paths is on the list.*

### 3 Proof of Theorem 1

In this section we prove Theorem 1. Let us first try to clarify our proof plan. It is based on a sequence of transformations whose goal is to simplify the random variables until a case analysis becomes manageable. Known arguments (the *reduce support* operation that will be explained in our proof) show that we may assume that every random variable by itself is “simple” in the sense that it has small support. Hence it is reasonable to expect that if the number of random variables is a small constant, then the theorem (if true) can be proven by a “brute force” case analysis. For example, Samuels [6] mentions 25 cases that are to be considered if one is to prove Conjecture 2 for the case  $n = 4$ . In our proof we describe a merge operation that replaces two random variables by one random variable. However we do not perform this merge operation until

the number of random variables becomes small, because this operation might create random variables whose mean is larger than 1. After rescaling the random variables, this corresponds to decreasing  $\delta$  in the statement of Theorem 1 to be arbitrarily small, and drives the bound  $\delta/(1+\delta)$  to 0. Instead, we perform the merge operation until a step in which many random variables may still remain, but they have nice properties. Namely, all random variables except for perhaps one have roughly the same mean, and moreover, have small “surplus” (to be defined in the proof). Thereafter, case analysis becomes possible, through the use of Proposition 8 that allows us to analyse many such random variables as if they were just one random variable. It is interesting to note that even though our proof plan does not seem to allow us to prove tight bounds (due to the fact that we generate random variables with mean above 1), it does in fact provide optimal bounds when  $\delta$  is small ( $\delta < 1/12$ ). We shall comment more on our proof in Section 3.2, and now we present the proof itself.

**Proof:** Fix  $n$ ,  $\delta$ , and arbitrary nonnegative random variables  $X_1, \dots, X_n$  with means at most 1. We prove that inequality (1) holds. We assume without loss of generality that the support of every random variable is composed of a finite set of values. (This is a standard argument, but we sketch it for completeness. Any value larger than  $\mu + \delta$  in the support of a random variable can be lowered to  $\mu + \delta$ , without increasing the probability that  $X < \mu + \delta$ . Thereafter, any continuous random variable can be approximated by a discrete random variable with the same mean and whose support includes only multiples of  $\epsilon$ , where  $\epsilon$  is chosen to be much smaller than  $\delta/n$ . For these new random variables,  $X'_1, \dots, X'_n$ , the event  $X' < \mu + \delta'$  where  $\delta' = \delta - \epsilon n$  implies that for the original variables,  $X < \mu + \delta$ . By making  $\epsilon$  arbitrarily small, we can make  $\delta'$  arbitrarily close to  $\delta$ .)

Our proof of inequality (1) consists of a sequence of transformations on the variables  $X_i$ . We may view these transformations as occurring in discrete time steps, and in our notation, superscripts will denote time steps. Hence, after time step  $t$ , random variables are denoted by  $X_i^t$ , their sum by  $X^t$  and the expectation of  $X^t$  by  $\mu^t$ . For  $t = 0$ , we have the original random variables. All our transformations will have the property that for every  $t \geq 0$ ,

$$Pr[X^{t+1} < \mu^{t+1} + \delta] \leq Pr[X^t < \mu^t + \delta]. \quad (2)$$

Some properties of the random variables may change by the transformations. In particular, the reduce support transformation (to be defined shortly) when applied to two random variables that were originally identically distributed might transform them to new random variables that are not identically distributed. Moreover, the merge transformation might generate random variables whose mean is larger than 1, even though all original random variables have mean at most 1. We now describe the transformations.

*Remove constant.* This transformation is applied whenever there is a random variable  $X_i^t$  that is constant, that is,  $Pr[X_i^t = \mu_i^t] = 1$ . Such a random variable is removed, and  $\mu^{t+1} = \mu^t - \mu_i^t$ . Clearly, *remove constant* satisfies inequality (2).

*Reduce support.* This transformation is applied to every random variable whose support has at least three values, and replaces it with a new random

variable with the same mean, and whose support includes at most two values from the original support.

**Lemma 4** *Let  $X_i^t$  be a random variable whose support includes at least three values. Then  $X_i^t$  can be replaced by a new variable  $X_i^{t+1}$  with  $\mu_i^{t+1} = \mu_i^t$ , and whose support includes only two values from the original support of  $X_i^t$ . This can be done in a way that satisfies inequality (2).*

**Proof:** Let  $\{v_1, \dots, v_k\}$  be the support of  $X_i^t$ , and for  $1 \leq j \leq k$ , let  $q_j$  denote the conditional probability of the event  $[X^t < \mu^t + \delta^t]$ , conditioned on the event  $[X_i^t = v_j]$ . For  $X_i^{t+1}$  and for  $1 \leq j \leq k$ , we wish to select  $p_j = Pr[X_i^{t+1} = v_j]$ , under the restrictions that the mean of  $X_i^{t+1}$  is the same as the mean of  $X_i^t$ , and that inequality (2) is satisfied. This can be expressed by the following linear program over the variables  $p_j$ :

**Minimize**  $\sum_{j=1}^k q_j p_j$   
subject to:

- $\sum_{j=1}^k p_j = 1$
- $\sum_{j=1}^k p_j v_j = \mu_i^t$
- $p_j \geq 0$ , for every  $j$ .

The above linear program is feasible (as the probabilities associated with the original  $X_i^t$  satisfy the constraints). By the theory of linear programming, there is a *basic* optimal solution in which at most two  $p_j$  are nonzero.  $\square$

We remark that Lemma 4 has several alternative proofs, and in general does not require  $X_t$  to have finite support. A similar lemma (with a functional analytic proof) is used in [6].

*Align with 0.* This transformation is applied to every random variable whose support has two values and these values are greater than 0 (say  $X_i^t$  has value  $v_1$  with probability  $p$  and  $v_2$  with probability  $(1-p)$ , with  $0 < v_1 < v_2$ ), and replaces it by a random variable  $X_i^{t+1}$  that has value  $v_1 - v_1 = 0$  with probability  $p$ , and has value  $v_2 - v_1$  with probability  $(1-p)$ . Hence  $\mu_i^{t+1} = \mu_i^t - v_1$ , and  $\mu^{t+1} = \mu^t - v_1$ . Clearly, *align with 0* satisfies inequality (2).

*Merge.* This transformation takes the two random variables with smallest mean (say  $X_i^t$  and  $X_j^t$ ), and replaces them by a new random variable in three steps. First, replace  $X_i^t$  and  $X_j^t$  by a new random variable that is distributed like their sum  $X_i^t + X_j^t$ . Then apply *reduce support* to this new random variable. Finally, apply *align with 0* or *remove constant* to the new random variable (if applicable).

It is easy to see that the transformation *merge* satisfies inequality (2).

The sequence of transformations that we perform is partitioned into two stages. We now describe the first stage.

**Stage 1:**

1. Whenever possible, apply *remove constant*.

2. Apply *reduce support* until all random variables have support of size at most two. (Different variables may have different support.)
3. Apply *align with 0* to all variables.
4. Apply *merge* until either the number of random variables is reduced to one, or all random variables have mean at least  $1/2$  (whichever happens first).

Stage 1 must end because with each application of *merge*, the number of random variables decreases. Let  $t$  denote the step after which stage 1 ends, and let  $X_1^t, \dots, X_{n'}^t$  be the random variables that remain. We assume that they are sorted in order of decreasing  $\mu_i^t$ . Their number  $n'$  may be smaller than  $n$ , because some of the transformations remove random variables. These are not arbitrary random variables, as each of them has a support of two values, one of which is 0, and the stopping condition for the merge transformations has been reached. For a random variable  $X_i^t$  as above, let  $\mu_i^t$  denote its mean,  $\{0, v_i^t\}$  its support, and let  $s_i^t = v_i^t - \mu_i^t$  denote its *surplus*. Let  $s^t = \sum_{i=1}^{n'} s_i^t$  denote the total surplus.

**Proposition 5** *If the total surplus satisfies  $s^t < \delta$ , then  $\Pr[X^t \geq \mu^t + \delta] = 0$ .*

**Proof:**  $X^t$  is maximized when  $X_i^t = v_i^t$  for all  $i$ . In this case

$$X^t = \sum_{i=1}^{n'} (\mu_i^t + s_i^t) = \mu^t + s^t < \mu^t + \delta.$$

□

Hence we may assume without loss of generality that  $s^t \geq \delta$ .

**Lemma 6** *If stage 1 ended with a random variable with mean below  $1/2$ , then  $\Pr[X^t < \mu^t] \geq \delta/(1/2 + \delta)$ .*

**Proof:** In this case, exactly one random variable remains. Let  $X_1^t$  be the random variable left, with mean  $\mu_1^t < 1/2$  and support  $\{0, v_1 = \mu_1^t + s^t\}$ . Note that the event  $X_1^t = 0$  implies  $X^t < \mu^t$ . Now  $\Pr[X_1^t = 0] = s^t/(\mu_1^t + s^t) \geq \delta/(1/2 + \delta)$ , because  $s^t \geq \delta$  and  $\mu_1^t < 1/2$ . □

Observe that Lemma 6 offers a conclusion that is even stronger than that required by Theorem 1, as is illustrated by the following sequence of inequalities.

$$\Pr[X < \mu + \delta] \geq \Pr[X^t < \mu^t + \delta] \geq \Pr[X^t < \mu^t] \geq \frac{\delta}{1/2 + \delta} > \frac{\delta}{1 + \delta}.$$

Hence we may also assume that stage 1 ended with all random variables having mean at least  $1/2$ . The following property will be used in this case.

**Proposition 7** *If stage 1 ended with all random variables having mean at least  $1/2$ , then  $\mu_1^t/2 \leq \mu_{n'}^t \leq \mu_1^t < 3/2$ .*



**Proof:** Recall that the random variables are assumed to be sorted with  $\mu_1^t$  being the largest mean and  $\mu_{n'}^t$  being the smallest mean.

If no random variable has mean greater than 1, then we are done. Hence consider the first time that a random variable with mean greater than 1 is created. This happens by merging two random variables, say at time step  $r$  (shortly we will see that in fact it must hold that  $r = t - 1$ ), with the random variables being  $X_i^r$  and  $X_j^r$ . Let  $\mu_i^r \geq \mu_j^r$  be their means before the merge. By the definition of *merge*, no other variable had mean smaller than  $\mu_i^r$ . By the stopping rule for stage 1,  $\mu_j^r < 1/2$ . To get a variable with mean greater than 1, we must have  $\mu_i^r > 1/2$ . Note that stage 1 ends after the merge, because no variable with mean below 1/2 is left. Hence the new variable created becomes  $X_1^{r+1}$  with  $1 < \mu_1^{r+1} < 1 + 1/2 = 3/2$ . But as  $\mu_1^{r+1} \leq 2\mu_i^r$  and  $\mu_{n'}^{r+1} \geq \mu_i^r$ , it follows that  $\mu_{n'}^{r+1} \geq \mu_1^{r+1}/2$ .  $\square$

Let us pause at this point and explain what remains to be proved. All random variables can be assumed to be 2-valued, with one of the values being 0, and with all means  $\mu_i^t$  satisfying  $\mu_1^t/2 \leq \mu_i^t \leq \mu_1^t$ . Moreover, the total surplus  $s^t$  satisfies  $s^t \geq \delta$ . For random variables as above we in fact will bound  $Pr[X^t < \mu^t]$  rather than  $Pr[X^t < \mu^t + \delta]$ . Lemma 9 (its first part) and Lemma 10 will show that  $Pr[X^t < \mu^t] \geq \min[\delta/(\mu_1^t + \delta), 1/13]$ . This almost proves Theorem 1, except that it might happen that at the end of stage 1,  $\mu_1^t > 1$ . This possibility is handled in the second part of Lemma 9, by showing that one merge operation before the end of stage 1 we had  $Pr[X^{t-1} < \mu^{t-1} + \delta] \geq \delta/(1 + \delta)$ .

The following proposition is used several times in the proofs of Lemmas 9 and 10. It is most effective when  $s < \mu_n$ , and  $\mu_n$  is not much smaller than  $\mu_1$ .

**Proposition 8** *Let  $X_1, \dots, X_n$  be independent random variables with means  $\mu_1 \geq \dots \geq \mu_n$  and supports  $\{0, \mu_1 + s_1\}, \dots, \{0, \mu_n + s_n\}$ , and let  $X = \sum_{i=1}^n X_i$ ,  $\mu = \sum_{i=1}^n \mu_i$  and  $s = \sum_{i=1}^n s_i$ . Then*

$$Pr[X < \mu - \mu_n + s] \geq \frac{s}{\mu_1 + s}$$

**Proof:** It suffices that one random variable comes up zero to imply  $X < \mu + s - \mu_n$ . (The inequality is strict because only a variable with  $s_i > 0$  may come up 0.) Hence:

$$Pr[X \geq \mu + s - \mu_n] = \prod_{i=1}^n \frac{\mu_i}{\mu_i + s_i} \leq \prod_{i=1}^n \frac{\mu_1}{\mu_1 + s_i}.$$

Given that  $\sum_{i=1}^n s_i = s$  and that  $s_i \geq 0$ , the above product is maximized when  $s_1 = s$  and  $s_i = 0$  for all  $i > 1$ , giving  $\mu_1/(\mu_1 + s)$ . Hence  $Pr[X < \mu - \mu_n + s] \geq s/(\mu_1 + s)$ .  $\square$

The following lemma illustrates the desired outcome of stage 1.

**Lemma 9** *1. If stage 1 ended with all random variables having mean at least 1/2, and if  $s^t < \mu_{n'}^t$ , then*

$$Pr[X^t < \mu^t] \geq \frac{\delta}{\mu_1^t + \delta} \geq \frac{\delta}{3/2 + \delta}.$$

2. If in addition  $\delta \leq 1/12$  then either  $\Pr[X^t < \mu^t] < \delta/(1+\delta)$ , or one merge operation before the end of stage 1 it must have been the case that

$$\Pr[X^{t-1} < \mu^{t-1} + \delta] \geq \delta/(1+\delta).$$

**Remark:** The choice of  $\delta \leq 1/12$  in the second part of Lemma 9 is made because  $\delta/(1+\delta) = 1/13$  for  $\delta = 1/12$ . The limiting factor for improving beyond  $1/13$  is Lemma 10 rather than Lemma 9. For  $\delta > 1/12$  the second part of Lemma 9 simply implies that  $\Pr[X^{t-1} < \mu^{t-1} + \delta] \geq \Pr[X^{t-1} < \mu^{t-1} + 1/12] \geq 1/13$ .

**Proof:** The surplus  $s^t$  is smaller than the mean of any of the random variables. Using Proposition 8 we then have  $\Pr[X^t < \mu^t] \geq \frac{s^t}{\mu_1^t + s^t}$ . Using the assumption that  $s^t \geq \delta$  and the fact that  $\mu_1^t \leq 3/2$  (Proposition 7), we have that  $\Pr[X^t < \mu^t] \geq \delta/(3/2 + \delta)$ .

To prove the second part of the lemma, note that if it happens that  $\mu_1^t \leq 1$  then we have  $\Pr[X^t < \mu^t] \geq \delta/(1+\delta)$ . Hence we may assume that  $\mu_1^t > 1$ , implying in particular that  $X_1^t$  is the result of the last *merge* operation (see proof of Proposition 7). Let  $s' = s^t - s_1^t$  be the surplus of all variables except for  $X_1^t$ . If  $s' > 0$  (which implies  $\mu_2^t \geq 1/2 > 0$ ) then analysis as in the proof of the first part of the Lemma implies that

$$\Pr[X^t \geq \mu^t] \leq \frac{\mu_2^t}{\mu_2^t + s'} \leq \frac{1}{1 + s'}.$$

Hence if  $s' \geq \delta$ ,  $\Pr[X^t < \mu^t] \geq \delta/(1+\delta)$ . So we can assume that  $s' < \delta$  (including also the possibility that  $s' = 0$ ).

Let us backtrack the last merge operation. Hence instead of  $X_1^t$  we have two variables  $X_i^{t-1}$  and  $X_j^{t-1}$  that were merged to give  $X_1^t$ . Let their means be  $\mu_i^{t-1} \geq \mu_j^{t-1}$ , and their surpluses be  $s_i^{t-1}$  and  $s_j^{t-1}$ . Observe that necessarily  $\mu_j^{t-1} < 1/2$  (otherwise the merge operation would not have been performed), and then the assumption that  $\mu_1^t > 1$  implies that  $\mu_i^{t-1} > 1/2$ . As the total surplus of all random variables except for  $X_i^{t-1}$  and  $X_j^{t-1}$  is  $s' < \delta$ , we must have  $X_i^{t-1} + X_j^{t-1}$  come up larger than  $\mu_i^{t-1} + \mu_j^{t-1}$  for  $X^{t-1} \geq \mu^{t-1} + \delta$ . We consider now two cases.

**Case 1:**  $s_i^{t-1} > 2\delta$ . Then  $\Pr[X_i^{t-1} = 0] = \frac{s_i^{t-1}}{\mu_i^{t-1} + s_i^{t-1}} \geq \frac{2\delta}{\mu_i^{t-1} + 2\delta}$ . If  $X_i^{t-1} = 0$ , then in order to have  $X_i^{t-1} + X_j^{t-1} > \mu_i^{t-1} + \mu_j^{t-1}$  we must have  $X_j^{t-1} > \mu_i^{t-1} + \mu_j^{t-1}$ . But this happens with probability at most  $\frac{\mu_j^{t-1}}{\mu_i^{t-1} + \mu_j^{t-1}} \leq \frac{1/2}{\mu_i^{t-1} + 1/2}$ . Hence

$$\Pr[X^{t-1} < \mu^{t-1} + \delta] \geq \frac{2\delta}{\mu_i^{t-1} + 2\delta} \cdot \frac{\mu_i^{t-1}}{\mu_i^{t-1} + 1/2} \geq \frac{\delta}{1 + \delta}$$

where the last inequality holds for  $\delta \leq 1/2$  because  $1/2 \leq \mu_i^{t-1} \leq 1$ .

**Case 2:**  $s_i^{t-1} \leq 2\delta$ . Define  $s'' = s' + s_i^{t-1}$  as the surplus of all random variables except for  $s_j^{t-1}$ , and observe that  $s'' < 3\delta \leq 1/4$ , the last inequality holding for  $\delta \leq 1/12$ . We now consider several subcases.

- $s'' < \delta$  and  $s_j^{t-1} \geq 1/2$ . The fact that  $s'' < \delta$  implies that it suffices for  $X_j^{t-1}$  to come up 0 to ensure  $X^{t-1} < \mu^{t-1} + \delta$ . This happens with probability at least  $s_j^{t-1}/(\mu_j^{t-1} + s_j^{t-1}) \geq 1/2$ .
- $s'' < \delta$  and  $s_j^{t-1} < 1/2$ . The fact that all random variables except for  $X_j^{t-1}$  have mean at least  $1/2$  implies that it suffices for one random variable to come up 0 to ensure  $X^{t-1} < \mu^{t-1} + \delta$ . As necessarily  $s'' + s_j^{t-1} \geq \delta$  and  $\mu_k^{t-1} \leq 1$  for all  $k$ , this happens with probability at least  $\delta/(1 + \delta)$ , by Proposition 8.
- $\delta \leq s'' < \delta + \mu_j^{t-1}$  and  $s_j^{t-1} \geq \delta$ . It suffices for  $X_j^{t-1}$  to come up 0 to ensure  $X^{t-1} < \mu^{t-1} + \delta$ . This happens with probability at least  $\delta/(1/2 + \delta)$ .
- $\delta \leq s'' < \delta + \mu_j^{t-1}$  and  $s_j^{t-1} < \delta$ . It suffices for some random variable other than  $X_j^{t-1}$  to come up 0 to ensure  $X^{t-1} < \mu_{t-1} - 1/2 + s'' + \delta < \mu^{t-1} + \delta$ . (We used the facts that every random variable except for  $X_j^{t-1}$  has mean at least  $1/2$ , and that  $s'' \leq 1/4$ .) This happens with probability at least  $\delta/(1 + \delta)$ , by Proposition 8.
- $s'' \geq \delta + \mu_j^{t-1}$ . By Proposition 8, there is probability of at least  $\frac{\delta + \mu_j^{t-1}}{1 + \delta + \mu_j^{t-1}}$  for a random variable other than  $X_j^{t-1}$  to come up 0. Thereafter  $X_j^{t-1}$  must come up at least  $\mu_j^{t-1} + \delta + 1/2 - s'' \geq \mu_j^{t-1} + \delta + 1/4$  for  $X^{t-1} \geq \mu^{t-1} + \delta$  to hold. The probability of this is at most  $\mu_j^{t-1}/(\mu_j^{t-1} + 1/4 + \delta) \leq \mu_j^{t-1}/(\mu_j^{t-1} + 1/4)$ . Hence

$$Pr[X^{t-1} < \mu^{t-1} + \delta] \geq \frac{\delta + \mu_j^{t-1}}{1 + \delta + \mu_j^{t-1}} \cdot \frac{1/4}{\mu_j^{t-1} + 1/4} \geq \frac{\delta}{1 + \delta}$$

where the last inequality holds when  $4(\delta + \delta^2 + \mu_j^{t-1}\delta) \leq 1$ , which is true for our parameters of  $\mu_j^{t-1} \leq 1/2$  and  $\delta \leq 1/12$ .

□

Summarizing, Lemmas 6 and 9 prove Theorem 1 except for the case that stage 1 ended with  $\mu_1^t/2 \leq \mu_{n'}^t \leq \mu_1^t$  and  $s^t \geq \mu_{n'}^t$ . We now address this last case. For this case we shall not use the extra slackness offered by  $\delta$ , but rather we shall show the stronger inequality  $Pr[X^t < \mu^t] \geq 1/13$ . To prove this last inequality, we perform stage 2 of our sequence of transformations. It is composed of a modified form of the merge operations, that we call *modified merge*. The modification will allow us to deal with the event  $X^t < \mu^t$  rather than  $X^t < \mu^t + \delta$ . Recall that the *reduce support* operation was based on a linear program that minimized  $Pr[X^{t+1} < \mu^{t+1} + \delta]$  (via the definition of the  $q_j$  in the proof of Lemma 4). Modify the *reduce support* operation by modifying the objective function of the linear program to be  $Pr[X^{t+1} < \mu^{t+1}]$  (by making the respective change in the definition of  $q_j$ ). Use this *modified reduce support*

rather than the original *reduce support* as the second step of *modified merge*. Now *modified merge* maintains the inequality

$$\Pr[X^{t+1} < \mu^{t+1}] \leq \Pr[X^t < \mu^t]. \quad (3)$$

Note that an application of *modified merge* cannot increase the total surplus (this was true also for *merge*). This works in our favor, because as the proof of Lemma 9 demonstrates, it is easier to perform case analysis when the total surplus is small. However, an application of *modified merge* may result in a random variable whose mean is smaller than  $\mu_1^{t+1}/2$ . (For simplicity of notation, we assume that after every step the variables are renamed so as to keep  $\mu_1^{t+1}$  the largest mean.) This will complicate our case analysis. But note that even with repeated applications of *modified merge*, there will be at most one such random variable. Let us define  $s' = \sum s_i$ , where the sum is taken over all random variables whose mean is at least  $\mu_1/2$ . In particular, at the time when stage 1 ends we may assume that  $s' = s^t$ , by Proposition 7.

**Stage 2.** Apply *modified merge* (on the two random variables with currently lowest mean) until a step (that we shall denote by  $r$ ) after which either the number of remaining (nonconstant) random variables is one, or the condition  $s' \leq \alpha\mu_1^r$  has been reached, for some constant  $0 < \alpha < 1/2$  that will be determined later. Stage 2 must eventually end, because with each application of *modified merge* the number of random variables decreases.

**Lemma 10** *Let  $\alpha = 1/3$ , and let  $r$  denote the time step at which stage 2 ends. Then either  $\Pr[X^r < \mu^r] \geq 1/13$ , or one modified merge operation before stage 2 ends  $\Pr[X^{r-1} < \mu^{r-1}] \geq 1/13$ .*

**Proof:** The proof of Lemma 10 is based on a case analysis. Some of the details in the case analysis are included so as to get the explicit bound of  $1/13$ . Those readers who are just interested in verifying that the lemma holds for some universal constant (though perhaps much smaller than  $1/13$ ) may simply think of  $\alpha$  as some small constant (say  $\alpha = 1/10$ ) and  $\beta$  (to be introduced shortly) as a much smaller constant (say  $\beta = 1/100$ ), and read each case only up to the point where it becomes clear that under this setting of the parameters, the case in question gives a probability bounded away from 0.

**Lemma 11** *If stage 2 ends without the condition  $s' \leq \alpha\mu_1^r$  being reached, then  $\Pr[X^r < \mu^r] > \alpha/(1 + \alpha)$ .*

**Proof:** In this case we have only one nonconstant random variable,  $X_1^r$ , with support  $\{0, \mu_1^r + s'\}$ .

$$\Pr[X^r < \mu^r] = \Pr[X_1^r = 0] = \frac{s'}{\mu_1^r + s'} > \frac{\alpha}{1 + \alpha}.$$

□

**Lemma 12** *If stage 2 ends with  $\beta\mu_1^r \leq s' \leq \alpha\mu_1^r$ , where  $0 \leq \beta \leq \alpha$  is some constant that will be optimized later, then*

$$Pr[X^r < \mu^r] \geq \min \left[ \left( \frac{\alpha - 2\alpha^2}{1 + \alpha} \right), \left( \frac{\beta - 2\beta^2}{1 + \beta} \right) \right].$$

**Proof:** Consider first only the random variables with mean at least  $\mu_1^r/2$ , let  $X'$  be their sum and let  $\mu'$  be the expectation of  $X'$ . Over these random variables, the surplus is  $s' = \gamma\mu_1^r$ , with  $\beta \leq \gamma \leq \alpha < 1/2$ . By Proposition 8,

$$Pr[X' < \mu' - (1/2 - \gamma)\mu_1^r] \geq \frac{s'}{\mu_1^r + s'} = \frac{\gamma}{1 + \gamma}.$$

The event  $X' < \mu' - (1/2 - \gamma)\mu_1^r$  does not yet imply that  $X^r < \mu^r$ . There still might be one variable  $X_{n''}^r$  with  $\mu_{n''}^r < \mu_1^r/2$ . If  $X_{n''}^r$  turns out  $\mu_{n''}^r + s_{n''}^r$  and  $s_{n''}^r \geq (1/2 - \gamma)\mu_1^r$  then it still may hold that  $X^r \geq \mu^r$ .

Let us first assume that  $\mu_{n''}^r \leq s' = \gamma\mu_1^r$ . Then by Markov's inequality,

$$Pr[X_{n''}^r = \mu_{n''}^r + s_{n''}^r] = \frac{\mu_{n''}^r}{\mu_{n''}^r + s_{n''}^r} \leq \frac{\gamma\mu_1^r}{\gamma\mu_1^r + (1/2 - \gamma)\mu_1^r} = 2\gamma.$$

Hence

$$Pr[X^r < \mu^r] \geq \frac{\gamma}{1 + \gamma} \cdot (1 - 2\gamma) = \frac{\gamma - 2\gamma^2}{1 + \gamma}.$$

For  $0 < \beta \leq \gamma \leq \alpha < 1/2$ , the expression above is minimized when  $\gamma \in \{\alpha, \beta\}$ .

We are left with the case that  $\mu_{n''}^r > s'$ . But then we have

$$Pr[X^r < \mu^r] \geq Pr[X_{n''}^r = 0] = \frac{s_{n''}^r}{\mu_{n''}^r + s_{n''}^r} \geq \frac{1/2 - \gamma}{1 - \gamma}$$

where we have used the facts that  $\mu_{n''}^r < \mu_1^r/2$  and  $s_{n''}^r \geq (1/2 - \gamma)\mu_1^r$ . As  $\gamma \leq \alpha$ , we have that  $Pr[X^r < \mu^r] \geq (1/2 - \alpha)/(1 - \alpha)$ . But this probability is larger than  $(\alpha - 2\alpha^2)/(1 + \alpha)$  of the previous case, and hence can be ignored.  $\square$

**Lemma 13** *If stage 2 ends with  $s' < \beta\mu_1^r$ , and  $0 < \beta < \alpha/2$ , then one merge prior to the end of stage 2 it must have been the case that  $Pr[X^{r-1} < \mu^{r-1}]$  was at least the minimum of the following expressions:*

1.  $\frac{\alpha - \beta}{1/2 + \alpha - \beta} \cdot \frac{1/2 - \beta}{1 - \beta}$
2.  $\frac{\alpha - 3\beta/2}{1 + \alpha - 3\beta/2} \cdot \frac{1 - 3\beta/2}{3/2 - 3\beta/2}$
3.  $\frac{\alpha - 2\beta}{1 + \alpha - 2\beta}$
4.  $\left( \frac{1/2 - \beta}{1 - \beta} \right)^2$
5.  $\frac{1/2 - 3\beta/2}{3/2 - 3\beta/2} \cdot \frac{1 - 3\beta/2}{3/2 - 3\beta/2}$

The proof of Lemma 13 involves a detailed case analysis and appears in Section 3.1.

Summing up, Lemma 12 and 13 imply that after stage 2, either  $Pr[X^r < \mu^r]$  or  $Pr[X^{r-1} < \mu^{r-1}]$  is at least the smallest of the following quantities (where  $0 < \beta < \alpha/2 < 1/4$ ):

- $\frac{\alpha-2\alpha^2}{1+\alpha}$
- $\frac{\beta-2\beta^2}{1+\beta}$
- $\frac{\alpha-\beta}{1/2+\alpha-\beta} \cdot \frac{1/2-\beta}{1-\beta}$
- $\frac{\alpha-3\beta/2}{1+\alpha-3\beta/2} \cdot \frac{1-3\beta/2}{3/2-3\beta/2}$
- $\frac{\alpha-2\beta}{1+\alpha-2\beta}$
- $\left(\frac{1/2-\beta}{1-\beta}\right)^2$
- $\frac{1/2-3\beta/2}{3/2-3\beta/2} \cdot \frac{1-3\beta/2}{3/2-3\beta/2}$

Choosing (suboptimally)  $\alpha = 1/3$  and  $\beta = 1/8$  gives at least  $1/13$  in all cases.

This completes the proof of Lemma 10.  $\square$

We can now summarize the proof of Theorem 1. We have the original random variables for which we wish to prove  $Pr[X < \mu + \delta] \geq \min[\delta/(1 + \delta), 1/13]$ . The proof proceeds in two stages. In stage 1 we apply a sequence of transformations maintaining the inequality (2) until a step  $t$  in which all random variables have mean at least  $1/2$  (or only one random variable remains, which is an easy case to handle). Then Lemma 9 implies the theorem whenever the total surplus  $s^t$  is small ( $s^t < \mu_n^t$ ). To handle the case that the surplus is large, we perform a sequence of transformations in stage 2, this time maintaining the inequality (3), until a step  $r$  in which the surplus of all variables (except for the variable with smallest mean) is no longer large compared to the maximum mean (namely,  $s' \leq \alpha\mu_1^r$ ). Then Lemma 10 implies that also in this case (corresponding to the case that  $s^t$  was greater than  $\mu_n^t$ ) the theorem holds. Note that overall, depending on which case is considered, Lemmas 9 and 10 do not prove inequality (1) on the original random variables, but on transformed random variables that appear after one of the steps  $t-1$ ,  $t$ ,  $r-1$  or  $r$ . But this implies inequality (1) on the original random variables, because of the inequalities (2) and (3).  $\square$

### 3.1 Remaining analysis for stage 2

We prove here Lemma 13, whose proof was the only part omitted from the proof of Theorem 1.

**Proof:** Consider the last two random variables to have been merged, say  $X_i^{r-1}$  and  $X_j^{r-1}$ , with means  $\mu_i^{r-1} \geq \mu_j^{r-1}$ , and let  $\mu_1^{r-1}$  be the largest mean

at time  $r - 1$ . After the *modified merge* of  $X_i^{r-1}$  and  $X_j^{r-1}$ , the largest mean  $\mu_1^r$  may still have been  $\mu_1^{r-1}$ , but it could also be as high as  $\mu_i^{r-1} + \mu_j^{r-1}$ , if this happens to be higher than  $\mu_1^{r-1}$ . In fact,  $\mu_1^r$  may also be lower than  $\mu_1^{r-1}$ , if only one variable is left at the end of stage 2, and this variable underwent an *align with 0* operation. However, in this case the bounds that we get for  $X^{r-1} < \mu^{r-1}$  are much stronger than what we get otherwise (details omitted), so we shall ignore this case.

We analyse the situation one merge operation before the end of stage 2. Note that we know that at that time,  $s' \geq \alpha\mu_1^{r-1}$ , because otherwise stage 2 would have ended earlier. Likewise, the sum  $\sum s_i^{r-1}$  taken over all variables except  $X_i^{r-1}$  and  $X_j^{r-1}$  is at most  $\max[\beta\mu_1^{r-1}, \beta(\mu_i^{r-1} + \mu_j^{r-1})]$ , because otherwise we could not have had  $s' < \beta\mu_1^r$  at the end of stage 2. We consider now two cases.

**Case 1:**  $\mu_j^{r-1} < \mu_1^{r-1}/2$ . Hence  $X_j^{r-1}$  did not contribute to  $s'$  at time step  $r - 1$ . Note that  $\mu_i^{r-1} \geq \mu_1^{r-1}/2$  and hence  $X_i^{r-1}$  did contribute to  $s'$  at time step  $r - 1$ . It follows that  $s_i^{r-1} \geq \alpha\mu_1^{r-1} - \beta\mu_1^r$ . (This last expression is positive, because Lemma 13 assumes that  $\beta < \alpha/2$ , and  $\mu^r \leq 2\mu^{r-1}$ .) Hence,

$$Pr[X_i^{r-1} = 0] \geq \frac{\alpha\mu_1^{r-1} - \beta\mu_1^r}{\mu_i^{r-1} + \alpha\mu_1^{r-1} - \beta\mu_1^r}$$

If  $X_i^{r-1} = 0$  then in order to have  $X^{r-1} \geq \mu^{r-1}$ ,  $X_j^{r-1}$  must contribute at least  $\mu_j^{r-1} + \mu_i^{r-1} - \beta\mu_1^r$  to  $X^{r-1}$ . (The expression  $\mu_i^{r-1} - \beta\mu_1^r$  is positive because  $\beta < 1/4$  and  $\mu_i^{r-1} \geq \mu_1^{r-1}/2 \geq \mu_1^r/4$ .) This may happen with probability at most  $\mu_j^{r-1}/(\mu_j^{r-1} + \mu_i^{r-1} - \beta\mu_1^r)$ . We then have

$$Pr[X^{r-1} < \mu^{r-1}] \geq \frac{\alpha\mu_1^{r-1} - \beta\mu_1^r}{\mu_i^{r-1} + \alpha\mu_1^{r-1} - \beta\mu_1^r} \cdot \frac{\mu_i^{r-1} - \beta\mu_1^r}{\mu_j^{r-1} + \mu_i^{r-1} - \beta\mu_1^r}$$

The above expression is minimized when  $\mu_j^{r-1}$  is maximized, namely, when  $\mu_j^{r-1} = \mu_1^{r-1}/2$ . As  $\mu_i^{r-1} \geq \mu_1^{r-1}/2$ , it follows that  $\mu_i^{r-1} + \mu_j^{r-1} \geq \mu_1^{r-1}$ . The expression above is minimized when  $\mu_1^r$  is maximized, namely,  $\mu_1^r = \mu_j^{r-1} + \mu_i^{r-1}$ . Normalising  $\mu_1^{r-1}$  to 1, and using the notation  $\mu_i$  to denote  $\mu_i^{r-1}/\mu_1^{r-1}$ , we have after some rearrangements

$$Pr[X^{r-1} < \mu^{r-1}] \geq \frac{\alpha - \beta/2 - \beta\mu_i}{(1 - \beta)\mu_i + \alpha - \beta/2} \cdot \frac{(1 - \beta)\mu_i - \beta/2}{(1 - \beta)\mu_i + 1/2 - \beta/2}$$

The expression above is defined for all  $\mu_i \geq 0$ . It equals 0 for  $\mu_i \in \{\beta/2(1 - \beta), (\alpha - \beta/2)/\beta\}$  and is positive in between. Moreover, there are at most two points where the derivative with respect to  $\mu_i$  of this expression vanishes (as it is a ratio of two non-proportional quadratics), and for  $\beta < 2\alpha/3$  the expression is positive in the allowed range of  $1/2 \leq \mu_i \leq 1$ . It follows that the expression is minimized when  $\mu_i \in \{1/2, 1\}$ , giving

$$Pr[X^{r-1} < \mu^{r-1}] \geq$$

$$\min \left[ \frac{\alpha - \beta}{1/2 + \alpha - \beta} \cdot \frac{1/2 - \beta}{1 - \beta}, \frac{\alpha - 3\beta/2}{1 + \alpha - 3\beta/2} \cdot \frac{1 - 3\beta/2}{3/2 - 3\beta/2} \right]$$

This gives items 1 and 2 in the statement of Lemma 13.

**Case 2:**  $\mu_j^{r-1} \geq \mu_1^{r-1}/2$ . Hence both  $s_i^{r-1}$  and  $s_j^{r-1}$  did contribute to  $s'$  (before the last merge), and moreover,  $\mu_i^{r-1} + \mu_j^{r-1} \geq \mu_1^{r-1}$ . This together with the inequality  $\mu_1^r \leq \max[\mu_1^{r-1}, \mu_i^{r-1} + \mu_j^{r-1}]$  imply that we may use  $\mu_i^{r-1} + \mu_j^{r-1}$  as an upper bound on  $\mu_1^r$ . To simplify notation and without loss of generality we may assume that  $\mu_1^{r-1} = 1$ , and then  $1/2 \leq \mu_j^{r-1} \leq \mu_i^{r-1} \leq 1$ . We have that  $s_i^{r-1} + s_j^{r-1} > \alpha - \beta(\mu_i^{r-1} + \mu_j^{r-1}) \geq \alpha - 2\beta$ . Recall (from the paragraph prior to case 1) the sum  $\sum s_i^{r-1}$  taken over all variables except  $X_i^{r-1}$  and  $X_j^{r-1}$  is at most  $\beta(\mu_i^{r-1} + \mu_j^{r-1})$ . Hence if  $X_i^{r-1} + X_j^{r-1} < \mu_i^{r-1} + \mu_j^{r-1} - \beta(\mu_i^{r-1} + \mu_j^{r-1})$ , then necessarily  $X^{r-1} < \mu^{r-1}$ . We let  $B$  denote the event  $[X_i^{r-1} + X_j^{r-1} < (1 - \beta)(\mu_i^{r-1} + \mu_j^{r-1})]$ , and perform a subcase analysis for  $Pr[B]$ . The subcases are partitioned according to which of  $X_i^{r-1}$  and  $X_j^{r-1}$  (or both, or either one) needs to come up 0 in order for  $B$  to hold.

1. It suffices that either  $X_i^{r-1} = 0$  or  $X_j^{r-1} = 0$  for  $B$  to hold. In this case, using  $\mu_i^{r-1} \leq 1$  and  $s_i^{r-1} + s_j^{r-1} \geq \alpha - 2\beta$ , Proposition 8 implies that

$$Pr[B] \geq \frac{\alpha - 2\beta}{1 + \alpha - 2\beta}$$

This gives item 3 in the statement of Lemma 13.

2.  $B$  holds iff  $X_i^{r-1} = 0$ . In this subcase necessarily  $\mu_i^{r-1} + s_i^{r-1} \geq (1 - \beta)(\mu_i^{r-1} + \mu_j^{r-1})$ . Using the fact that  $\mu_i^{r-1} \leq 1$  and  $\mu_j^{r-1} \geq 1/2$  we have that  $s_i^{r-1} \geq \mu_j^{r-1} - \beta\mu_i^{r-1} - \beta\mu_j^{r-1} \geq 1/2 - 3\beta/2$  and therefore

$$Pr[B] = Pr[X_i^{r-1} = 0] = \frac{s_i^{r-1}}{\mu_i^{r-1} + s_i^{r-1}} \geq \frac{1/2 - 3\beta/2}{3/2(1 - \beta)}$$

This subcase is dominated by the subcase above and hence can be ignored.

3.  $B$  holds iff  $X_j^{r-1} = 0$ . This subcase is analogous to and dominated by the subcase above, and can be ignored.
4.  $B$  holds only if both  $X_i^{r-1} = 0$  and  $X_j^{r-1} = 0$ . Then necessarily  $\mu_i^{r-1} + s_i^{r-1} \geq (1 - \beta)(\mu_i^{r-1} + \mu_j^{r-1})$  and  $\mu_j^{r-1} + s_j^{r-1} \geq (1 - \beta)(\mu_i^{r-1} + \mu_j^{r-1})$ . We have

$$Pr[B] \geq \left( \frac{(1 - \beta)(\mu_i^{r-1} + \mu_j^{r-1}) - \mu_i^{r-1}}{(1 - \beta)(\mu_i^{r-1} + \mu_j^{r-1})} \right) \cdot \left( \frac{(1 - \beta)(\mu_i^{r-1} + \mu_j^{r-1}) - \mu_j^{r-1}}{(1 - \beta)(\mu_i^{r-1} + \mu_j^{r-1})} \right)$$

For fixed  $\mu_i^{r-1} + \mu_j^{r-1}$  this expression is minimized when  $\mu_i^{r-1} - \mu_j^{r-1}$  is maximized. Hence either  $\mu_i^{r-1} = 1$  or  $\mu_j^{r-1} = 1/2$ . Thereafter, it can be



verified that the expression is minimized when the other mean is either maximized or minimized, giving us three possible local minimum points,  $\mu_i^{r-1}, \mu_j^{r-1} \in \{1/2, 1\}$ ,  $\mu_j^{r-1} \leq \mu_i^{r-1}$ . Two of these give identical values (the cases that  $\mu_i^{r-1} = \mu_j^{r-1}$ ), hence we obtain

$$Pr[B] \geq \min \left[ \left( \frac{1/2 - \beta}{1 - \beta} \right)^2, \left( \frac{1/2 - 3\beta/2}{3/2 - 3\beta/2} \cdot \frac{1 - 3\beta/2}{3/2 - 3\beta/2} \right) \right]$$

This gives items 4 and 5 in the statement of Lemma 13.

□

### 3.2 Some comments

It is straightforward to modify inequality (1) so that there is no formal requirement that the random variables are nonnegative, or that their mean is bounded by 1. Let  $w$  be the maximum over all random variables  $X_1, \dots, X_n$  of the respective  $\mu_i - l_i$ , where  $l_i$  is the lowest value in the support of  $X_i$ . Then

$$Pr[X \leq \mu + \delta w] \geq \min[\delta/(1 + \delta), 1/13] \quad (4)$$

The constant 1/13 in Theorem 1 is not best possible, and can be improved with more detailed case analysis. We suspect that the true constant should be  $1/e$ . Presumably, the way to prove a tight result is to find a sequence of transformations on the random variables that does not increase  $Pr[X < \mu + \delta]$ , and that gradually brings them to the conjectured worst case for  $[X < \mu + \delta]$ . The sequence of transformations performed in our proof of Theorem 1 manages to achieve this only when  $\delta \leq 1/12$  (or some other constant not far from  $1/12$ ). However, it fails to characterize the worst case for the perhaps more interesting  $\delta = 1$ . The idea in the proof is to transform the random variables into a situation where a case analysis becomes manageable, at the possible cost of giving up the tightness of the bound. The main principles used are reducing the support of every random variable to two values, getting all random variables (perhaps except one) to have roughly the same mean, reducing the surplus to be of order of magnitude comparable to this mean, and extracting from arbitrarily many random variables a single event of interest, as done in Proposition 8. It should be clear to the reader that more detailed case analysis would provide tighter results. But let us point out some limitations that relate to Lemma 10. As long as one chooses  $\alpha$  not larger than  $\mu_1/2$  (and in fact, not larger than  $3\mu_1/2$ ), and analyses only the situation at the end of stage 2 or one step earlier, one cannot obtain a bound better than  $Pr[X < \mu] \geq 2/9$ . For example, assume that during stage 2 we are left with three variables, each with support  $\{0, \mu\}$  and mean  $\mu/3$ . At this point,  $Pr[X < \mu] = (2/3)^3 = 8/27 < 1/e$ . After a merge operation, this probability decreases further to  $(1/3) \cdot (2/3) = 2/9$ . One merge operation later, stage 2 ends. Hence to get (nearly) tight results using the current approach, one may need to modify the definition of stage 2, and perform much more extensive (possibly computer assisted) case analysis.

## 4 Proof of Theorem 2

The reader is assumed to be familiar with elementary methods in probability (such as the use of Markov's inequality, Chebyshev's inequality, Chernoff bounds). If needed, see details in [1], for example.

We query at random  $t$  vertices, and obtain their degrees. Let  $d_i$  be the degree returned by the  $i$ th query. Basically, our estimator for  $d$  will be  $d^* = \frac{1}{t} \sum_{i=1}^t d_i$ . In section 4.3 we shall modify this estimator so as to improve its quality. For simplicity of the analysis, we assume that sampling is done with replacement (the same vertex might be queried more than once). This is insignificant when  $t$  is small (e.g.,  $t \leq \sqrt{n}$ ), though note that for large values of  $t$  (and in particular, when  $t = n$ ) sampling without replacement gives better estimates than sampling with replacement.

Note that the expectation of our estimate satisfies

$$E[d^*] = d \tag{5}$$

Hence the estimator is *unbiased*. In deviations from the expectation, we will analyse separately the events  $d^* > d$  and  $d^* < d$ , or rather,  $d^* < d/2$ .

### 4.1 The estimate is not too high

Here we shall use Theorem 1. As an immediate consequence of this theorem (taking  $\delta = 1$ , and using the fact that the degree of a sampled vertex is a non-negative random variable with expectation  $d$ ) we have the following corollary.

**Corollary 14** *There is some universal constant  $\alpha > 0$  such that for every graph with average degree  $d$ , by querying  $t$  random vertices (with replacement) for their degree, the average  $d^*$  satisfies  $\Pr[d^* \leq (1 + 1/t)d] \geq \alpha$ .*

We can take  $\alpha = 1/13$  in Corollary 14, and we conjecture that the Corollary is also true with  $\alpha = 1/e$ .

### 4.2 The estimate is not too low

We assume that the average degree in the graph is at least  $d_0$ . Our sampling algorithm queries  $t = k\sqrt{n/d_0}$  vertices at random and reports the sum of the degrees. Here  $k$  is a parameter that will later be chosen to be of order  $1/\epsilon$ .

Let  $X_i$  be the random variable that denotes the degree of the  $i$ th query, and let  $X = \sum_{i=1}^t X_i$ . Then  $E[X] = t \cdot d$ . The following lemma shows that the typical value of  $X$  is not much smaller than  $E[X]/2$ .

**Lemma 15** *For arbitrary  $\delta > 0$  (that will later be fixed to  $50\sqrt{2}/\alpha$ , where  $\alpha$  is as in Corollary 14), with probability at least  $1 - 4\sqrt{2}/\delta - 2^{-\Omega(\delta)}$ ,*

$$X \geq \frac{E[X]}{2} \left(1 - \frac{\delta}{k}\right)$$

**Proof:** Essentially, the proof of the lemma is based on Chebyshev's inequality. To apply Chebyshev's inequality directly, one would need the variance of  $X$  to be small compared to  $(E[X])^2$ . Unfortunately, vertices of very high degree may cause the variance to exceed  $(E[X])^2$ . To overcome this problem we observe (and will soon show formally) that in every graph, the vertices of very high degree contain at most slightly more than half the endpoints of the edges. The contribution to  $X$  of vertices whose degree is not very high is concentrated around its mean, because for them the variance is small. This explains why the value of  $X$  is likely to be not much smaller than  $E[X]/2$ . We now present our proof of the lemma, based on the above principles.

Partition the set of vertices of  $G$  into two sets,  $H$  (for high) and  $L$  (for low). For a constant  $c$  (independent of  $n, d, k$ ) that will be determined later, the set  $H$  contains the  $c\sqrt{nd}/k$  vertices of highest degree (breaking ties arbitrarily). The set  $L$  contains the other vertices. Every edge has two endpoints. Let us partition the endpoints of edges into the following four sets:

- $E_{H,L}$  (the endpoints in  $H$  of edges between  $H$  and  $L$ )
- $E_{L,H}$  (the endpoints in  $L$  of edges between  $H$  and  $L$ )
- $E_{H,H}$  (the endpoints in  $H$  of edges between  $H$  and  $H$ )
- $E_{L,L}$  (the endpoints in  $L$  of edges between  $L$  and  $L$ )

Observe that  $|E_{H,H}| \leq |H|^2 = c^2nd/k^2$ . It will be the case that  $c$  is a universal constant whereas  $k \geq \Omega(1/\epsilon)$ , and hence  $|E_{H,H}| = O(\epsilon^2nd)$ . Moreover, we allow an error of  $\epsilon \cdot nd$  in our estimation of  $nd$ . Hence,  $E_{H,H}$  has only a low order effect on the accuracy of the estimation. So as to simplify notation and the presentation, we shall simply assume that  $|E_{H,H}| = 0$ . We shall not give a rigorous proof that this assumption has only a low order effect on our analysis, but merely note here that formalists may redo the analysis without assuming that  $|E_{H,H}| = 0$ , and at worst this will effect some constants that are eventually hidden by the  $O$  notation.

Let  $m_1 = |E_{H,L}|$ ,  $m_2 = |E_{L,H}|$  and  $m_3 = |E_{L,L}|$ . Hence  $m_1 + m_2 + m_3 = dn$ . Note that  $m_2 = m_1$ , because  $|E_{L,H}| = |E_{H,L}|$ . Let us break the random variable  $X$  into the sum of three random variables  $X = Y_1 + Y_2 + Y_3$ , according to the contribution to  $X$  from  $m_1$ ,  $m_2$  and  $m_3$  respectively. Let  $h$  denote the minimum degree of a vertex in  $H$ .

**Proposition 16** *With probability  $1 - 2^{-\Omega(c)}$ ,*

$$Y_1 \geq ch/2$$

**Proof:** The expected number of vertices queried from  $H$  is  $t|H|/n = k\sqrt{n/d_0}$ .  $c\sqrt{d_0n}/kn = c$ . With probability  $1 - 2^{-\Omega(c)}$ , the actual number of vertices queried from  $H$  is at least  $c/2$ . Each such vertex contributes at least  $h$  to  $Y_1$ .  $\square$

**Proposition 17** A vertex in  $L$  can cover at most  $|H| = c\sqrt{nd_0}/k$  endpoints in  $E_{L,H}$ .

**Proof:** For every endpoint in  $E_{L,H}$  covered by a vertex in  $L$ , the other endpoint of the respective edge is in  $H$ . As the original graph is a simple graph with no parallel edges, the proof follows.  $\square$

**Proposition 18** For  $\lambda > 0$ , with probability at least  $1 - 1/\lambda^2$ ,

$$Y_2 \geq E[Y_2] - \lambda\sqrt{cdn/2}$$

**Proof:** The variance of  $Y_2$  is maximized if the endpoints of  $E_{L,H}$  are concentrated on  $m_2/|H|$  vertices (each covering  $|H|$  endpoints). Hence:

$$\text{var}[Y_2] \leq |H|^2 \frac{m_2}{n|H|} t = cm_2 \leq cdn/2$$

The proof now follows from Chebyshev's inequality.  $\square$

**Proposition 19** For  $\lambda > 0$ , with probability at least  $1 - 1/\lambda^2$ ,

$$Y_3 \geq E[Y_3] - \lambda\sqrt{\frac{hkm_3}{\sqrt{d_0n}}}$$

**Proof:** The maximum degree of any vertex in  $L$  is  $h$ . Hence the graph induced by the edges  $E_{L,L}$  also has maximum degree at most  $h$ . Thus

$$\text{var}[Y_3] \leq h^2 \frac{m_3}{h} \frac{t}{n} = hm_3k/\sqrt{d_0n}$$

The proof now follows from Chebyshev's inequality.  $\square$

**Proposition 20** With probability at least  $1 - 2/\lambda^2 - 2^{-\Omega(c)}$ ,

$$X \geq \frac{E[X]}{2} + \frac{ch}{2} - \lambda\sqrt{\frac{cdn}{2}} - \lambda\sqrt{\frac{hkm_3}{\sqrt{d_0n}}} + \frac{km_3}{2\sqrt{d_0n}}$$

**Proof:**  $X = Y_1 + Y_2 + Y_3$ . By propositions 16,18 and 19 we have that with probability at least  $1 - 2/\lambda^2 - 2^{-\Omega(c)}$ ,

$$X \geq E[Y_2] + E[Y_3] + \frac{ch}{2} - \lambda\sqrt{\frac{cdn}{2}} - \lambda\sqrt{\frac{hkm_3}{\sqrt{d_0n}}}$$

As  $E[Y_1] = E[Y_2]$ , we have that  $E[X]/2 = E[Y_2] + E[Y_3]/2$ . Using  $E[Y_3] = m_3t/n = km_3/\sqrt{d_0n}$  the proof follows.  $\square$

Fix  $c = 4\lambda^2$ . Then

$$\frac{ch}{2} \cdot \frac{km_3}{2\sqrt{d_0n}} \geq \left( \lambda\sqrt{\frac{hkm_3}{\sqrt{d_0n}}} \right)^2$$

implying

$$\frac{ch}{2} - \lambda \sqrt{\frac{hkm_3}{d_0 n}} + \frac{km_3}{2\sqrt{d_0 n}} \geq 0.$$

Hence for  $c = 4\lambda^2$ , the inequality in Proposition 20 can be replaced by

$$X \geq \frac{E[X]}{2} - \lambda \sqrt{\frac{cdn}{2}}$$

The term  $\lambda \sqrt{\frac{cdn}{2}} = \lambda^2 \sqrt{2dn}$  is at most  $\frac{E[X]}{2} 2\lambda^2 \sqrt{2}/k$ , because  $E[X] = dt = k\sqrt{nd}\sqrt{d/d_0}$ . Renaming  $2\sqrt{2}\lambda^2$  by  $\delta$ , Lemma 15 is proved.  $\square$

### 4.3 Combining the upper and lower bound

Let us set  $k = 3\delta/\epsilon$  (where  $\epsilon$  is taken from Theorem 2), and hence from Lemma 15 we have that with probability at least  $1 - 4\sqrt{2}/\delta - 2^{-\Omega(\delta)}$ ,

$$X \geq \frac{E[X]}{2}(1 - \epsilon/3)$$

By Corollary 14, we have that with probability at least  $\alpha$ ,  $X \leq E[X](1 + 1/t)$ . The ratio between the upper bound on  $X$  and the lower bound on  $X$  is  $2(1 + 1/t)/(1 - \epsilon/3) \leq 2 + \epsilon$ . This last inequality holds when  $\epsilon$  is sufficiently small (which implies Theorem 2 also for larger values of  $\epsilon$ ), and  $t$  is sufficiently large compared to  $1/\epsilon$  (which is true in our context because  $t$  is a parameter that grows with the number of vertices  $n$ , and the  $O$  notation in the statement of Theorem 2 implies that it suffices to prove the theorem when  $n$  is sufficiently large).

An *unbiased estimate* consists of taking  $t$  samples and returning their sum  $X$ . Perform  $2/\alpha$  independent unbiased estimates for  $X$ , where  $\alpha$  is taken to be as in the discussion following Corollary 14. Our estimation procedure returns  $X_{\min}$ , the minimum of these estimates. (Equivalently, we set  $d^* = X_{\min}/t$ .)

$$\Pr[X_{\min} \leq E[X](1 + \frac{1}{t})] \geq 1 - (1 - \alpha)^{2/\alpha} \geq 1 - \frac{1}{e^2} \geq \frac{5}{6}$$

$$\Pr[X_{\min} \geq \frac{E[X]}{2}(1 - \epsilon/3)] \geq 1 - \frac{2}{\alpha} \left( \frac{4\sqrt{2}}{\delta} + 2^{-\Omega(\delta)} \right) \geq 5/6$$

where the last inequality uses  $\delta = 50\sqrt{2}/\alpha$ . This gives  $k = 3\delta/\epsilon < 220/\alpha\epsilon$ . The total number of queries used in our estimation procedure is  $2t/\alpha$ . This gives:

**Corollary 21** *For some universal constant  $\beta$ , using*

$$\beta \frac{\sqrt{n/d_0}}{\epsilon}$$

*queries, one can estimate the average degree  $d$  of an  $n$  node graph within a ratio of  $(2 + \epsilon)$ , provided that  $d > d_0$ .*

**Proof:** Setting  $\beta = (\frac{2}{\alpha})(\frac{220}{\alpha}) = 440/\alpha^2$ , we perform  $2/\alpha$  unbiased estimates, each with  $t = 220\sqrt{n/d_0}\alpha^{-1}\epsilon^{-1}$  queries, and take the minimum of the estimations that they give.  $\square$

Let us note here the role of Corollary 14. It allows us to substitute a universal constant for  $\alpha$  (which is shown to be at least  $1/13$  in Theorem 1, though we conjecture that  $1/e$  also works). If not for Theorem 1, we could have used Markov's inequality in a proof of a modified Corollary 14, showing (for example) that  $\Pr[d^* \leq (1 + \epsilon/3)d] \geq \epsilon/4$ . This would have been equivalent to replacing  $\alpha$  in the proof of Corollary 21 by  $\epsilon/4$ , which would require the number of queries used by the estimation procedure to be  $\beta\sqrt{n/d_0}\epsilon^{-3}$  (for some constant  $\beta$ ). This a factor of  $\epsilon^{-2}$  worse than the bounds that we get through the use of Theorem 1.

#### 4.4 Optimality of sample size

The sample size in Corollary 21 is essentially best possible, as the following proposition shows.

**Proposition 22** *For every (reasonable)  $n, d, \epsilon$ , one can construct a graph  $G_1$  with  $(1 + \epsilon)nd$  edges and a graph  $G_2$  with  $dn/2$  edges, such that  $\Omega(\frac{1}{\epsilon}\sqrt{n/d})$  vertices need to be queried in order to have probability above  $2/3$  of distinguishing between them.*

**Proof:** Graph  $G_1$  has a set  $A$  of  $\epsilon\sqrt{nd}$  vertices of degree  $(1 + \epsilon)\sqrt{nd}/\epsilon$ , and a set  $B$  of  $(1 + \epsilon)\sqrt{nd}/\epsilon$  vertices of degree  $\epsilon\sqrt{nd}$  (e.g., arranged as a complete bipartite subgraph between  $A$  and  $B$ ). The other vertices have degree 0. Graph  $G_2$  has a set  $C$  of  $\sqrt{nd}/\epsilon$  vertices of degree  $\epsilon\sqrt{nd}$ .

We sketch the proof of why  $\Omega(\sqrt{\frac{n}{d}}\epsilon^{-1})$  queries are necessary. Assume that the number of queries is  $\sqrt{\frac{n}{d}}\epsilon^{-1}$ . Then there is constant probability that no vertex from  $A$  is queried, and the expected number of vertices queried from  $B$  is  $\epsilon^{-2} + \epsilon^{-1}$ . The expected vertices queried from  $C$  is  $\epsilon^{-2}$ . As the standard deviation is of order  $\sqrt{\epsilon^{-2}} = \epsilon^{-1}$ , there is constant probability that  $G_1$  and  $G_2$  will be confused.  $\square$

The optimality of the sample size was proved under the assumption that the only information used by the estimation algorithm is the degree of the queried vertices. More generally, one may think of randomized estimation algorithms that make use of additional information. For example, when querying a vertex of positive degree, the next vertex to query may be chosen at random from the list of neighbors of the current vertex. The use of a more general class of random estimation algorithms may allow either quicker or more accurate estimation of the average degree in a graph. See [5], for example. However, let us explain here some the the advantages of “degree only” sampling, advantages that might be lost by other estimation algorithms.

1. All queries can be made in parallel, which in some contexts results in saving time.

2. Sampling can be done anonymously. The estimation algorithm need not know the identity of queried vertices, nor the identity of their neighbors. Privacy issues may sometimes require that this be the case. For example, vertices of a graph may represent persons in some community, and an edge may represent some sort of interaction that took place between the respective persons. Persons may be willing to fill an anonymous questionnaire stating with how many different persons they had interaction (namely, their degree), but may not be willing to disclose with whom they had interaction.
3. In Section 5 there are several different graphs  $G_e$  defined on the same set of vertices, and in a single degree query one gets the degrees of the respective vertex in all graphs simultaneously. In order to efficiently estimate the average degree in all graphs, it is useful to have an estimation algorithm for which the choice of which vertex to query does not depend on the graph in question.

## 5 Quickly estimating the load on a network

We have seen how to estimate the average degree in a graph using a relatively small number of degree-queries. Graph problems are often abstractions of other more concrete problems. As an example (which motivated this study), consider the following problem motivated and studied in [3].

The input is a connected network  $G$  with  $n$  vertices and  $m$  edges (namely, a graph). Between every two vertices there is a shortest path (a path that crosses the smallest number of edges). We assume here that shortest paths are unique, a point that we shall return to later. For an integer parameter  $c$  (that may depend on  $n$ ), we wish to find all edges that are members of at least  $c$  shortest paths. In the terminology of [3], these edges are called “weakest links”, apparently because these are the edges where failure may cause the largest amount of damage to the performance of the network. Finding all weakest links can be done in time  $O(nm)$  using an algorithm for all pairs shortest paths. The goal in [3] is to do better. They propose a randomized algorithm that with high probability, has the following guarantee:

- **Finds weakest links.** It outputs all edges that belong to at least  $c$  shortest paths.
- **Avoids false alarms.** It does not output any edge that is a member of less than  $(1 - \epsilon)c$  shortest paths.

The running time of the algorithm in [3] is  $O(\frac{mn^2 \log n}{c\epsilon^2})$ , which is better than that of all pairs shortest paths when  $c \gg n \log n$ . The basic idea in this algorithm is to choose  $k \simeq \frac{n^2 \log n}{c\epsilon^2}$  pairs of vertices at random, and for each pair to perform a shortest path computation (taking  $O(m)$  operations per-pair). Using the collection of  $k$  shortest paths that are found, one estimates in how many shortest paths each edge participates.

Here we present a faster algorithm for finding the weakest links. It is based on two observations. One is that the cost of performing single source shortest path computations (namely, that of finding the shortest paths from one vertex to all other vertices) is  $O(m)$ , similar to that of finding the shortest path between one pair of vertices. The other observation is that the estimation problem that this gives rise to can be cast as that of estimating the average degree in a graph, or more precisely, in  $m$  different graphs simultaneously. The improved running time comes at a cost of a somewhat weaker guarantee in terms of false alarms.

- **Avoids false alarms.** The algorithm does not output any edge that is a member of less than  $(1/2 - \epsilon)c$  shortest paths.

As in [3] we assume that shortest paths are unique. This requires a convention for breaking ties between paths of equal length. We shall use the same convention that is proposed in [3], namely, to take the lexicographically first such path. This convention assumes that vertices are labelled and that there is a total ordering on the labels. For example, the labels can be 1 to  $n$ . A path can be viewed as a sequence of vertices in a natural way. Hence a path is a sequence of labels. In fact, two sequences correspond to the same path, depending on which of its two endpoints is considered to be the head of the path, and which is considered to be the tail. The name of the path is taken to be the lexicographically smaller of the two. Given two different paths that connect the same pair of vertices, if they are of equal length we use the convention that the one with the lexicographically smaller name is considered to be shorter.

**Proposition 23** *Under the tie breaking convention specified above, there is an  $O(m)$ -time algorithm that does the following. Given a connected graph  $G$  with  $n$  vertices and  $m$  edges and an arbitrary vertex  $v$ , it simultaneously counts for every edge  $e$ , for how many vertices  $u$  does edge  $e$  participate in the shortest path connecting  $u$  and  $v$ .*

**Proof:** We assume a model of computation in which algorithms such as single source shortest path take  $O(m)$  time. In particular, some basic operations (such as comparison between two  $O(\log n)$ -bit words) take unit time.

Given a starting vertex  $v$ , the distances to all other vertices in  $G$  can be computed in  $O(m)$  time using breadth first search (BFS). The BFS tree rooted at  $v$  gives also shortest paths from  $v$  to all vertices. It is quite straightforward to also count for each edge in the BFS tree (starting from edges furthest from the root and moving towards the root) in how many shortest paths (starting from  $v$ ) it participates. The counting requires only  $O(n)$  operations, as there are only  $n - 1$  edges in the BFS tree.

In general, several different BFS trees can be constructed starting at the same vertex  $v$ , because a vertex at distance  $i$  from  $v$  may have more than one neighbor at distance  $i - 1$  from  $v$ . We shall need to construct two such trees. For both trees, we may scan the vertices of the graph in the following order, starting at  $v$ : within a level of the BFS tree, vertices are scanned in the order under which they were first discovered, and every vertex scans its neighbors in order of



increasing labels. The *forward tree* rooted at  $v$  (gives the lexicographically first shortest paths when  $v$  is the first vertex of the path) is constructed using the following rule: for every vertex discovered at level  $i$  keep a pointer to its level  $i - 1$  neighbor that was first to be discovered (according to the scanning order described above). The *backward tree* rooted at  $v$  (gives the lexicographically first shortest paths when  $v$  is the last vertex of the path) is constructed using the following rule: for every vertex discovered at level  $i$  keep a pointer to its level  $i - 1$  neighbor of smallest label. Both the forward tree and the backward tree can be constructed in  $O(m)$  time.

Given both the forward tree and the backward tree for a vertex  $v$ , and using the convention that for vertices with a label smaller than  $v$  one uses the backward tree and for vertices with a label larger than  $v$  one uses the forward tree, one can simultaneously count in  $O(n)$  time how many shortest paths with an endpoint at  $v$  pass through every edge. (Note that this count is 0 for all but at most  $2n - 2$  edges of the two BFS trees.)  $\square$

We now consider  $m$  different graphs, one for every edge  $e$ . We denote the graph that we associate with edge  $e$  by  $G_e$ . The vertices of  $G_e$  are the vertices of  $G$ . Two vertices are connected by an edge in  $G_e$  iff  $e$  is on their unique shortest path in  $G$ . It follows that edge  $e$  is on  $c$  shortest paths in  $G$  iff the average degree in  $G_e$  is at least  $2c/n$ . Hence to find all weakest links, it suffices to find (or estimate) the average degrees of all graphs  $G_e$ . We shall now combine two facts.

1. By Theorem 2,  $O(\sqrt{n/d_0}/\epsilon)$  degree queries suffice in order to estimate the average degree in a graph with average degree at least  $d_0$ . To make the probability of error in this estimation below  $1/n^2$ , one can repeat the estimation procedure  $O(\log n)$  times, and take the median of the estimations. We shall set  $d_0 = (1 - \epsilon)c/n$ .
2. For any vertex  $v$ , Proposition 23 implies that in time  $O(m)$  one can simultaneously obtain the degree of  $v$  in all graphs  $G_e$ .

Hence using  $k = O(\frac{\log n \sqrt{n/(c/n)}}{\epsilon}) = O(\frac{n \log n}{\epsilon \sqrt{c}})$  single source shortest path computations one can with high probability simultaneously estimate the average degree in all graphs  $G_e$ , and by this find all weakest links (edges that are on more than  $c$  shortest paths) while avoiding any false alarms (by edges that are on less than  $(1/2 - \epsilon)c$  shortest paths). This proves Theorem 3.

Theorem 3 offers a saving of roughly  $n/\sqrt{c}$  in the running time compared to the running time of  $O(\frac{mn^2 \log n}{c\epsilon^2})$  in [3]. (Note however that  $\epsilon$  has different meanings in the two bounds. Hence the saving comes at the cost of allowing more false alarms.)

## Acknowledgements

I thank Johan Hastad, Michael Langberg, Eran Ofek, Gideon Schechtman, Benjy Weiss and Avi Wigderson for their interest and involvement in various

stages of this work, Assaf Naor for bringing to my attention the work of Samuels, and the anonymous referees for their detailed and useful comments.

## References

- [1] N. Alon and J. Spencer. *The Probabilistic Method*. Wiley-Interscience, 2000.
- [2] Z. W. Birnbaum, J. Raymond, H. S. Zuckerman. “A generalization of Tshebyshev’s inequality to two dimensions.” *Ann. Math. Statist.* Vol 18, 70–79, 1947.
- [3] N. Devanur, R. Lipton, N. Vishnoi. “Who’s the weakest link?” *Second Symposium on Stochastic Algorithms, Foundations and Applications*, SAGA 2003.
- [4] L. Dubins and L. Savage. *Inequalities for Stochastic Processes. (How to gamble if you must.)* Dover Publications, New York, 1976.
- [5] O Goldreich and D. Ron. “On estimating the average degree of a graph.” *Electronic Colloquim on Computational Complexity (ECCC)*, TR04-13, 2004.
- [6] S. M. Samuels. “On a Chebyshev-type inequality for sums of independent random variables.” *Ann. Math. Statist.* Vol 37, 248–259, 1966.
- [7] S. M. Samuels. “The Markov inequality for sums of independent random variables.” *Ann. Math. Statist.* Vol 40, 1980–1984, 1969.
- [8] A. Siegel. “Median bounds and their application.” *Journal of Algorithms*, 38(1), 184–236, 2001.