

# Randomized Approximation Algorithms for Set Multicover Problems with Applications to Reverse Engineering of Protein and Gene Networks\*

Piotr Berman<sup>†</sup>

Bhaskar DasGupta<sup>‡</sup>

Eduardo Sontag<sup>§</sup>

August 10, 2006

## Abstract

In this paper we investigate the computational complexity of a combinatorial problem that arises in the reverse engineering of protein and gene networks. Our contributions are as follows:

- We abstract a combinatorial version of the problem and observe that this is “equivalent” to the set multicover problem when the “coverage” factor  $k$  is a function of the number of elements  $n$  of the universe. An important special case for our application is the case in which  $k = n - 1$ .
- We observe that the standard greedy algorithm produces an approximation ratio of  $\Omega(\log n)$  even if  $k$  is “large” *i.e.*  $k = n - c$  for some constant  $c > 0$ .
- Let  $1 < a < n$  denote the maximum number of elements in any given set in our set multicover problem. Then, we show that a non-trivial analysis of a simple randomized polynomial-time approximation algorithm for this problem yields an expected approximation ratio  $\mathbf{E}[r(a, k)]$  that is an increasing function of  $a/k$ . The behavior of  $\mathbf{E}[r(a, k)]$  is roughly as follows: it is about  $\ln(a/k)$  when  $a/k$  is at least about  $e^2 \approx 7.39$ , and for smaller values of  $a/k$  it decreases towards 1 as a linear function of  $\sqrt{a/k}$  with  $\lim_{a/k \rightarrow 0} \mathbf{E}[r(a, k)] = 1$ . Our randomized algorithm is a cascade of a deterministic and a randomized rounding step parameterized by a quantity  $\beta$  followed by a greedy solution for the remaining problem. We also comment about the impossibility of a significantly faster convergence of  $\mathbf{E}[r(a, k)]$  towards 1 for any polynomial-time approximation algorithm.

**Keywords:** Set multicover, randomized approximation algorithms, reverse engineering, biological networks.

## 1 Introduction

Let  $[x, y]$  be the set  $\{x, x + 1, x + 2, \dots, y\}$  for integers  $x$  and  $y$ . The set multicover problem is a well-known combinatorial problem that can be defined as follows.

---

\*A preliminary version of these results with slightly weaker bounds appeared in 7th International Workshop on Approximation Algorithms for Combinatorial Optimization Problems, LNCS 3122, K. Jansen, S. Khanna, J. D. P. Rolim and D. Ron (editors), Springer Verlag, pp. 39-50, August 2004.

<sup>†</sup>Department of Computer Science and Engineering, Pennsylvania State University, University Park, PA 16802. Email: [berman@cse.psu.edu](mailto:berman@cse.psu.edu). Supported by NSF grant CCR-O208821.

<sup>‡</sup>Department of Computer Science, University of Illinois at Chicago, Chicago, IL 60607-7053. Email: [dasgupta@cs.uic.edu](mailto:dasgupta@cs.uic.edu). Supported in part by NSF grants CCR-0206795, CCR-0208749 and a CAREER grant IIS-0346973.

<sup>§</sup>Department of Mathematics, Rutgers University, New Brunswick, NJ 08903. Email: [sontag@hilbert.rutgers.edu](mailto:sontag@hilbert.rutgers.edu). Supported in part by NSF grant CCR-0206789.

**Problem name:**  $\mathbf{SC}_k$ .

**Instance**  $\langle n, m, k \rangle$ : An universe  $U = [1, n]$ , sets  $S_1, S_2, \dots, S_m \subseteq U$  with  $\cup_{j=1}^m S_j = U$  and a “coverage factor” (positive integer)  $k$ .

**Valid Solutions:** A subset of indices  $I \subseteq [1, m]$  such that, for every element  $x \in U$ ,  $|j \in I : x \in S_j| \geq k$ .

**Objective:** *Minimize*  $|I|$ .

$\mathbf{SC}_1$  is simply called the Set Cover problem and denoted by  $\mathbf{SC}$ ; we will denote an instance of  $\mathbf{SC}$  simply by  $\langle n, m \rangle$  instead of  $\langle n, m, 1 \rangle$ .

Both  $\mathbf{SC}$  and  $\mathbf{SC}_k$  are already well-known in the realm of design and analysis of combinatorial algorithms (e.g., see [18]). Let  $3 \leq a < n$  denote the maximum number of elements in any set, i.e.,  $a = \max_{i \in [1, m]} |S_i|$ . We summarize some of the known relevant results for them below.

### Fact 1<sup>1</sup>

(a) [6] *Assuming*  $\mathbf{NP} \not\subseteq \mathbf{DTIME}(n^{\log \log n})$ , instances  $\langle n, m \rangle$  of the  $\mathbf{SC}$  problem cannot be approximated to within a factor of  $(1 - \varepsilon) \ln n$  for any constant  $0 < \varepsilon < 1$  in polynomial time.

(b) [18] *An instance*  $\langle n, m, k \rangle$  of the  $\mathbf{SC}_k$  problem can be  $(1 + \ln a)$ -approximated in  $O(nmk)$  time by a simple greedy heuristic that, at every step, selects a new set that covers the maximum number of those elements that has not been covered at least  $k$  times yet. It is also possible to design randomized approximation algorithms with similar expected approximation ratios.

## 1.1 Summary of Results

The combinatorial problems investigated in this paper that arise out of reverse engineering of gene and protein networks can be shown to be equivalent to  $\mathbf{SC}_k$  when  $k$  is a function of  $n$ . One case that is of significant interest is when  $k$  is “large”, i.e.,  $k = n - c$  for some constant  $c > 0$ , but the case of non-constant  $c$  is also interesting (cf. Questions (Q1) and (Q2) in Section 2). Our contributions in this paper are as follows:

- In Section 2 we discuss the combinatorial problems (Questions (Q1) and (Q2)) with their biological motivations that are of relevance to the reverse engineering of protein and gene networks. We then observe, in Section 2.3, using a standard duality that these problems are indeed equivalent to  $\mathbf{SC}_k$  for appropriate values of  $k$ .
- In Lemma 2 in Section 3.1, we observe that the standard greedy algorithm  $\mathbf{SC}_k$  produces an approximation ratio of  $\Omega(\log n)$  even if  $k$  is “large”, i.e.  $k = n - c$  for some constant  $c > 0$ .
- Let  $1 < a < n$  denotes the maximum number of elements in any given set in our set multi-cover problem. In Theorem 3 in Section 3.2, we show that a non-trivial analysis of a simple randomized polynomial-time approximation algorithm for this problem yields an expected approximation ratio  $\mathbf{E}[r(a, k)]$  that is an increasing function of  $a/k$ . The behavior of  $\mathbf{E}[r(a, k)]$  is “roughly” as follows: it is about  $\ln(a/k)$  when  $a/k$  is at least about  $e^2 \approx 7.39$ , and for smaller values of  $a/k$  it decreases towards 1 as a linear function of  $\sqrt{a/k}$  with  $\lim_{a/k \rightarrow 0} \mathbf{E}[r(a, k)] = 1$ .

---

<sup>1</sup>A slightly weaker lower bound under the more standard complexity-theoretic assumption of  $\mathbf{P} \neq \mathbf{NP}$  was obtained by Raz and Safra [13] who showed that there is a constant  $c$  such that it is NP-hard to approximate instances  $\langle n, m \rangle$  of the  $\mathbf{SC}$  problem to within a factor of  $c \ln n$ .

More precisely,  $\mathbf{E}[r(\mathbf{a}, k)]$  is at most<sup>2</sup>

$$\begin{aligned}
& 1 + \ln \mathbf{a}, && \text{if } k = 1 \\
& (1 + e^{-(k-1)/5}) \ln(\mathbf{a}/(k-1)), && \text{if } \mathbf{a}/(k-1) \geq e^2 \approx 7.39 \text{ and } k > 1 \\
& \min\{2 + 2 \cdot e^{-(k-1)/5}, 2 + (e^{-2} + e^{-9/8}) \cdot \frac{\mathbf{a}}{k}\} \\
& \approx \min\{2 + 2 \cdot e^{-(k-1)/5}, 2 + 0.46 \cdot \frac{\mathbf{a}}{k}\} && \text{if } \frac{1}{4} < \mathbf{a}/(k-1) < e^2 \text{ and } k > 1 \\
& 1 + 2\sqrt{\frac{\mathbf{a}}{k}} && \text{if } \mathbf{a}/(k-1) \leq \frac{1}{4} \text{ and } k > 1
\end{aligned}$$

## 1.2 Summary of Analysis Techniques

- To prove Lemma 2, we generalize the approach in Johnson’s paper [8]. A straightforward replication of the sets will not work because of the dependence of  $k$  on  $n$ , but allowing the “misleading” sets to be somewhat larger than the “correct” sets allows a similar approach to go through at the expense of a diminished constant.
- Our randomized algorithm in Theorem 3 is a cascade of a deterministic and a randomized rounding step parameterized by a quantity  $\beta$  followed by a greedy solution for the remaining problem.
- Our analysis of the randomized algorithm in Theorem 3 uses an amortized analysis of the interaction between the deterministic and randomized rounding steps with the greedy step. For tight analysis, we found that the standard Chernoff bounds such as in [1, 3, 12, 18] were not always sufficient and hence we had to devise more appropriate bounds for certain parameter ranges.

## 1.3 Impossibility of Significantly Faster Convergence of $\mathbf{E}[r(\mathbf{a}, k)]$ Towards 1

It is certainly tempting to investigate the possibility of designing randomized or deterministic approximation algorithms for which  $\mathbf{E}[r(\mathbf{a}, k)]$  or  $r(\mathbf{a}, k)$  converges to 1 at a *significantly* faster rate as a function of  $\mathbf{a}/k$ . However, this may be difficult to achieve and, in particular,  $\mathbf{E}[r(\mathbf{a}, k)]$  or  $r(\mathbf{a}, k)$  cannot be  $1 + o(1)$  for  $\mathbf{a} \geq k$  since the set multicover problem is APX-hard for this case. To illustrate the last assertion, consider the special case of  $k = \mathbf{a} = n - 1$ . Then, the set multicover problem is still APX-hard as shown in the following. One could have  $n - 1$  sets of the form  $V \setminus \{i\}$  that cover every element, except one, exactly  $n - 2$  times (the last element is covered  $n - 1$  times). Moreover, we can have a family of sets of size exactly 3 that form an instance of the set cover problem restricted to  $\mathbf{a} = 3$ . This restricted problem is APX-hard, and a solution of size  $m$  for that instance gives solution of size  $n + m - 1$  for our instance. Because  $m \geq n/3$ , this is an approximation-preserving reduction. However, we will not investigate designing tight lower bounds further in this paper.

## 2 Motivations

In this section we define a computational problem that arises in the context of experimental design for reverse engineering of protein and gene networks. We will first pose the problem in linear algebra

<sup>2</sup>Note that, for  $k > 1$ , the bound on  $\mathbf{E}[r(\mathbf{a}, k)]$  is defined over three regions of values of  $\mathbf{a}/(k-1)$ , namely  $[\mathbf{a}, e^2)$ ,  $[e^2, \frac{1}{4})$  and  $[\frac{1}{4}, 0)$ . The boundaries between the regions can be shifted slightly by exact tedious calculations. We omit such straightforward but tedious exact calculations for simplicity.

terms, and then recast it as a combinatorial question. After that, we will discuss its motivations from systems biology. Finally, we will provide a precise definition of the combinatorial problems and point out its equivalence to the set multicover problem via a standard duality.

Our problem is described in terms of two matrices  $A \in \mathbb{R}^{n \times n}$  and  $B \in \mathbb{R}^{n \times m}$  such that:

- $A$  is *unknown*;
- $B$  is *initially unknown*, but each of its columns, denoted as  $B_1, B_2, \dots, B_m$ , can be retrieved with a *unit-cost query*;
- the columns of  $B$  are in *general position*, *i.e.*, each subset of  $\ell \leq n$  columns of  $B$  is *linearly independent*;
- the *zero structure* of the matrix  $C = AB = (c_{ij})$  is known, *i.e.*, a binary matrix  $C^0 = (c_{ij}^0) \in \{0, 1\}^{n \times m}$  is given, and it is known that  $c_{ij} = 0$  for each  $i, j$  for which  $c_{ij}^0 = 0$ .

The objective is to obtain as much information as possible about  $A$  (which, in the motivating application, describes regulatory interactions among genes and/or proteins), while performing “few” queries (each of which may represent the measuring of a complete pattern of gene expression, done under a different set of experimental conditions). For each query that we perform, we obtain a column  $B_i$ , and then the matrix  $C^0$  tells us that certain rows of  $A$  have zero inner product with  $B_i$ .

As a concrete example, let us take  $n = 3$ ,  $m = 5$ , and suppose that the known information is given by the matrix:

$$C_0 = \begin{bmatrix} 0 & 1 & 0 & 1 & 1 \\ 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 \end{bmatrix}$$

and the two unknown matrices are:

$$A = \begin{bmatrix} -1 & 1 & 3 \\ 2 & -1 & 4 \\ 0 & 0 & -1 \end{bmatrix}, \quad B = \begin{bmatrix} 4 & 3 & 37 & 1 & 10 \\ 4 & 5 & 52 & 2 & 16 \\ 0 & 0 & -5 & 0 & -1 \end{bmatrix}$$

(the matrix  $C_0$  has zero entries wherever  $AB$  has a zero entry). Considering the structure of  $C_0$ , we choose to perform four queries, corresponding to the four columns 1,3,4,5 of  $B$ , thus obtaining the following data:

$$\begin{bmatrix} 4 & 37 & 1 & 10 \\ 4 & 52 & 2 & 16 \\ 0 & -5 & 0 & -1 \end{bmatrix}. \quad (1)$$

What can we say about the unknown matrix  $A$ ? Let us first attempt to identify its first row, which we call  $A_1$ . The first row of the matrix  $C_0$  tells us that the vector  $A_1$  is orthogonal to the first and second columns of (1) (which are the same as the first and third columns of  $B$ ). This is the *only* information about  $A$  that we have available, and it is not enough information to uniquely determine  $A_1$ , because there is an entire line that is orthogonal to the plane spanned by these two columns. However, we can still find *some* nonzero vector in this line, and conclude that  $A_1$  is an unknown multiple of this vector. This nonzero vector may be obtained by simple linear algebra

manipulations. For example, we might add a linearly independent column to the two that we had, obtaining a matrix

$$B_1 = \begin{bmatrix} 4 & 37 & 0 \\ 4 & 52 & 0 \\ 0 & -5 & 1 \end{bmatrix},$$

then pick an arbitrary vector  $v$  whose first two entries are zero (to reflect the known orthogonality), let us say  $v = [0, 0, 1]$ , and finally solve  $A_1 B = v$ , thus estimating  $A_1$  as  $vB^{-1}$ :

$$\hat{A}_1 = [0, 0, 1] B^{-1} = [0, 0, 1] \begin{bmatrix} 13/15 & -37/60 & 0 \\ -1/15 & 1/15 & 0 \\ -1/3 & 1/3 & 1 \end{bmatrix} = [-1/3, 1/3, 1].$$

Notice that this differs from the unknown  $A_1$  only by a scaling. Similarly, we may employ the last two columns of (1) to estimate the second row  $A_2$  of  $A$ , again only up to a multiplication by a constant, and we may use the first and third columns of (1) (which are the same as the first and fourth columns of  $B$ ) to estimate the last row,  $A_3$ .

Notice that there are always intrinsic limits to what can be accomplished: if we multiply each row of  $A$  by some nonzero number, then the zero structure of  $C$  is unchanged. Thus, as in the example, the best that we can hope for is to identify the rows of  $A$  up to scalings (in abstract mathematical terms, as elements of the projective space  $\mathbb{P}^{n-1}$ ). To better understand these geometric constraints, let us reformulate the problem as follows. Let  $A_i$  denote the  $i^{\text{th}}$  row of  $A$ . Then the specification of  $C^0$  amounts to the specification of *orthogonality relations*  $A_i \cdot B_j = 0$  for each pair  $i, j$  for which  $c_{ij}^0 = 0$ . Suppose that we decide to query the columns of  $B$  indexed by  $J = \{j_1, \dots, j_\ell\}$ . Then, the information obtained about  $A$  may be summarized as  $A_i \in \mathcal{H}_{J,i}^\perp$ , where “ $\perp$ ” indicates *orthogonal complement*, and

$$\begin{aligned} \mathcal{H}_{J,i} &= \text{span} \{B_j, j \in J\}, \\ J_i &= \{j \mid j \in J \text{ and } c_{ij}^0 = 0\}. \end{aligned} \tag{2}$$

Suppose now that the set of indices of selected queries  $J$  has the property:

$$\text{each set } J_i, i = 1, \dots, n, \text{ has cardinality } \geq n - k, \tag{3}$$

for some given integer  $k$ . Then, because of the general position assumption, the space  $\mathcal{H}_{J,i}$  has dimension  $\geq n - k$ , and hence the space  $\mathcal{H}_{J,i}^\perp$  has dimension at most  $k$ .

### The case $k = 1$

The most desirable special case (and the one illustrated with the concrete example given above) is that in which  $k = 1$ . Then  $\dim \mathcal{H}_{J,i}^\perp \leq 1$ , hence each  $A_i$  is uniquely determined up to a scalar multiple, which is the best that could be theoretically achieved. Often, in fact, finding the sign pattern (such as “ $(+, +, -, 0, 0, -, \dots)$ ”) for each row of  $A$  is the main experimental goal (this would correspond, in our motivating application, to determining if the regulatory interactions affecting each given gene or protein are *inhibitory* or *catalytic*). Assuming that the degenerate case  $\mathcal{H}_{J,i}^\perp = \{0\}$  does not hold (which would determine  $A_i = 0$ ), once that an arbitrary nonzero element  $v$  in the line  $\mathcal{H}_{J,i}^\perp$  has been picked, there are only two sign patterns possible for  $A_i$  (the pattern of  $v$  and that of  $-v$ ). If, in addition, one knows at least one nonzero sign in  $A_i$ , then the sign structure of the whole row has been *uniquely* determined (in the motivating biological question, typically one

such sign is indeed known; for example, the diagonal elements  $a_{ii}$ , i.e. the  $i$ th element of each  $A_i$ , is known to be negative, as it represents a degradation rate). Thus, we will be interested in this question:

$$\text{find } J \text{ of minimal cardinality such that } |J_i| \geq n - 1, i = 1, \dots, n. \quad (\mathbf{Q1})$$

If queries have variable unit costs (different experiments have a different associated cost), this problem must be modified to that of minimizing a suitable linear combination of costs, instead of the number of queries.

### The general case $k > 1$

More generally, suppose that the queries that we performed satisfy (3), with  $k > 1$  but small  $k$ . It is not true anymore that there are only two possible sign patterns for any given  $A_i$ , but the number of possibilities is still very small. For simplicity, let us assume that we know that no entry of  $A_i$  is zero (if this is not the case, the number of possibilities may increase, but the argument is very similar). We wish to prove that the possible number of signs is much smaller than  $2^n$ . Indeed, suppose that the queries have been performed, and that we then calculate, based on the obtained  $B_j$ 's, a basis  $\{v_1, \dots, v_k\}$  of  $\mathcal{H}_{J,i}^\perp$  (assume  $\dim \mathcal{H}_{J,i}^\perp = k$ ; otherwise pick a smaller  $k$ ). Thus, the vector  $A_i$  is known to have the form  $\sum_{r=1}^k \lambda_r v_r$  for some (unknown) real numbers  $\lambda_1, \dots, \lambda_k$ . We may assume that  $\lambda_1 \neq 0$  (since, if  $A_i = \sum_{r=2}^k \lambda_r v_r$ , the vector  $\varepsilon v_1 + \sum_{r=2}^k \lambda_r v_r$ , with small enough  $\varepsilon$ , has the same sign pattern as  $A_i$ , and we are counting the possible sign patterns). If  $\lambda_1 > 0$ , we may divide by  $\lambda_1$  and simply count how many sign patterns there are when  $\lambda_1 = 1$ ; we then double this estimate to include the case  $\lambda_1 < 0$ . Let  $v_r = \text{col}(v_{1r}, \dots, v_{nr})$ , for each  $r = 1, \dots, k$ . Since no coordinate of  $A_i$  is zero, we know that  $A_i$  belongs to the set  $\mathcal{C} = \mathbb{R}^{k-1} \setminus (L_1 \cup \dots \cup L_n)$  where, for each  $1 \leq s \leq n$ ,  $L_s$  is the hyperplane in  $\mathbb{R}^{k-1}$  consisting of all those vectors  $(\lambda_2, \dots, \lambda_k)$  such that  $\sum_{r=2}^k \lambda_r v_{sr} = -v_{s1}$ . On each connected component of  $\mathcal{C}$ , signs patterns are constant. Thus the possible number of sign patterns is upper bounded by the maximum possible number of connected regions determined by  $n$  hyperplanes in dimension  $k - 1$ . A result of L. Schläfli (see [4, 14], and also [15] for a discussion, proof, and relations to Vapnik-Chervonenkis dimension) states that this number is bounded above by  $\Phi(n, k - 1)$ , provided that  $k - 1 \leq n$ , where  $\Phi(n, d)$  is the number of possible subsets of an  $n$ -element set with at most  $d$  elements, that is,

$$\Phi(n, d) = \sum_{i=0}^d \binom{n}{i} \leq 2 \frac{n^d}{d!} \leq \left(\frac{en}{d}\right)^d.$$

Doubling the estimate to include  $\lambda_1 < 0$ , we have the upper bound  $2\Phi(n, k - 1)$ . For example,  $\Phi(n, 0) = 1$ ,  $\Phi(n, 1) = n + 1$ , and  $\Phi(n, 2) = \frac{1}{2}(n^2 + n + 2)$ . Thus we have an estimate of 2 sign patterns when  $k = 1$  (as obtained earlier),  $2n + 2$  when  $k = 2$ ,  $n^2 + n + 2$  when  $k = 3$ , and so forth. In general, the number grows only polynomially in  $n$  (for fixed  $k$ ).

These considerations lead us to formulating the generalized problem, for each fixed  $k$ : *find  $J$  of minimal cardinality such that  $|J_i| \geq n - k$  for all  $i = 1, \dots, n$* . Recalling the definition (2) of  $J_i$ , we see that  $J_i = J \cap T_i$ , where  $T_i = \{j \mid c_{ij}^0 = 0\}$ . Thus, we can reformulate our question purely combinatorially, as a more general version of Question (Q1) as follows. Given sets

$$T_i \subseteq \{1, \dots, m\}, \quad i = 1, \dots, n.$$

and an integer  $k < n$ , the problem is:

$$\text{find } J \subseteq \{1, \dots, m\} \text{ of minimal cardinality such that } |J \cap T_i| \geq n - k, 1 \leq i \leq n. \quad (\mathbf{Q2})$$

For example, suppose that  $k = 1$ , and pick the matrix  $C^0 \in \{0, 1\}^{n \times n}$  in such a way that the columns of  $C^0$  are the binary vectors representing all the  $(n-1)$ -element subsets of  $\{1, \dots, n\}$  (so  $m = n$ ); in this case, the set  $J$  must equal  $\{1, \dots, m\}$  and hence has cardinality  $n$ . On the other hand, also with  $k = 1$ , if we pick the matrix  $C^0$  in such a way that the columns of  $C^0$  are the binary vectors representing all the 2-element subsets of  $\{1, \dots, n\}$  (so  $m = n(n-1)/2$ ), then  $J$  must again be the set of all columns (because, since there are only two zeros in each column, there can only be a total of  $2\ell$  zeros,  $\ell = |J|$ , in the submatrix indexed by  $J$ , but we also have that  $2\ell \geq n(n-1)$ , since each of the  $n$  rows must have  $\geq n-1$  zeros); thus in this case the minimal cardinality is  $n(n-1)/2$ .

## 2.1 Motivations from Systems Biology

This work was motivated by a central concern of contemporary cell biology, that of unraveling (or “reverse engineering”) the web of interactions among the components of complex protein and genetic regulatory networks. Notwithstanding the remarkable progress in genetics and molecular biology in the sequencing of the genomes of a number of species, the inference and quantification of interconnections in signaling and genetic networks that are critical to cell function is still a challenging practical and theoretical problem. High-throughput technologies allow the monitoring the expression levels of sets of genes, and the activity states of signaling proteins, providing snapshots of the transcriptional and signaling behavior of living cells. Statistical and machine learning techniques, such as clustering, are often used in order to group genes into co-expression patterns, but they are less able to explain functional interactions. An intrinsic difficulty in capturing such interactions in intact cells by traditional genetic experiments or pharmacological interventions is that any perturbation to a particular gene or signaling component may rapidly propagate throughout the network, causing global changes. The question thus arises of how to use the observed global changes to derive interactions between individual nodes.

This problem has generated an effort by many research groups whose goal is to infer mechanistic relationships underlying the observed behavior of complex molecular networks. We focus our attention here solely on one such approach, originally described in [10, 11], further elaborated upon in [2, 16], and reviewed in [5, 17]. In this approach, the architecture of the network is inferred on the basis of observed global responses (namely, the steady-state concentrations in changes in the phosphorylation states or activities of proteins, mRNA levels, or transcription rates) in response to experimental perturbations (representing the effect of hormones, growth factors, neurotransmitters, or of pharmacological interventions).

In the setup in [10, 11, 16], one assumes that the time evolution of a vector of state variables  $\mathbf{x}(t) = (x_1(t), \dots, x_n(t))$  is described by a system of differential equations:

$$\begin{aligned}\dot{x}_1 &= f_1(x_1, \dots, x_n, p_1, \dots, p_m) \\ \dot{x}_2 &= f_2(x_1, \dots, x_n, p_1, \dots, p_m) \\ &\vdots \\ \dot{x}_n &= f_n(x_1, \dots, x_n, p_1, \dots, p_m)\end{aligned}$$

(in vector form, “ $\dot{\mathbf{x}} = f(\mathbf{x}, \mathbf{p})$ ”, and the dot indicates time derivative), where  $\mathbf{p} = (p_1, \dots, p_m)$  is a vector of parameters, which can be manipulated but remain constant during any given experiment. The components  $x_i(t)$  of the state vector represent quantities that can be in principle measured, such as levels of activity of selected proteins or transcription rates of certain genes. The parameters  $p_i$  represent quantities that can be manipulated, perhaps indirectly, such as levels of hormones or of enzymes whose half-lives are long compared to the rate at which the variables evolve. A basic assumption (but see [16] for a time-dependent analysis) is that states converge to steady state

values, and these are the values used for network identification. There is a reference value  $\bar{p}$  of  $p$ , which represents “wild type” (that is, normal) conditions, and a corresponding steady state  $\bar{x}$ . Mathematically,  $f(\bar{x}, \bar{p}) = 0$ . We are interested in obtaining information about the Jacobian of the vector field  $f$  evaluated at  $(\bar{x}, \bar{p})$ , or at least about the signs of the derivatives  $\partial f_i / \partial x_j(\bar{x}, \bar{p})$ . For example, if  $\partial f_i / \partial x_j > 0$ , this means that  $x_j$  has a positive (catalytic) effect upon the rate of formation of  $x_i$ . The critical assumption, indeed the main point of [10, 11, 16], is that, while we may not know the form of  $f$ , we often do know that *certain parameters  $p_j$  do not directly affect certain variables  $x_i$* . This amounts to *a priori* biological knowledge of specificity of enzymes and similar data. In the current paper, this knowledge is summarized by the binary matrix  $C^0 = (c_{ij}^0) \in \{0, 1\}^{n \times m}$ , where “ $c_{ij}^0 = 0$ ” means that  $p_j$  does not appear in the equation for  $\dot{x}_i$ , that is,  $\partial f_i / \partial p_j \equiv 0$ .

The experimental protocol allows one to perturb any one of the parameters, let us say the  $k$ th one, while leaving the remaining ones constant. (A generalization, to allow for the simultaneous perturbation of more than one parameter, is of course possible.) For the perturbed vector  $p \approx \bar{p}$ , one then measures the resulting steady state vector  $x = \xi(p)$ . Experimentally, this may for instance mean that the concentration of a certain chemical represented by  $p_k$  is kept at a slightly altered level, compared to the default value  $\bar{p}_k$ ; then, the system is allowed to relax to steady state, after which the complete state  $x$  is measured, for example by means of a suitable biological reporting mechanism, such as a microarray used to measure the expression profile of the variables  $x_i$ . (Mathematically, we suppose that for each vector of parameters  $p$  in a neighborhood of  $\bar{p}$  there is a unique steady state  $\xi(p)$  of the system, where  $\xi$  is a differentiable function.)

For each of the possible  $m$  experiments, in which a given  $p_j$  is perturbed, we may estimate the  $n$  “sensitivities”

$$b_{ij} = \frac{\partial \xi_i}{\partial p_j}(\bar{p}) \approx \frac{1}{\bar{p}_j - p_j} (\xi_i(\bar{p} + p_j e_j) - \xi_i(\bar{p})) , \quad i = 1, \dots, n$$

(where  $e_j \in \mathbb{R}^m$  is the  $j$ th canonical basis vector). We let  $B$  denote the matrix consisting of the  $b_{ij}$ ’s. (See [10, 11] for a discussion of the fact that division by  $\bar{p}_j - p_j$ , which is undesirable numerically, is not in fact necessary.) Finally, we let  $A$  be the Jacobian matrix  $\partial f / \partial x$  and let  $C$  be the negative of the Jacobian matrix  $\partial f / \partial p$ . From  $f(\xi(p), p) \equiv 0$ , taking derivatives with respect to  $p$ , and using the chain rule, we get that  $C = AB$ . This brings us to the problem stated in this paper. (The general position assumption is reasonable, since we are dealing with experimental data.)

## 2.2 Combinatorial Formulation of Questions (Q1) and (Q2)

**Problem name:**  $\mathbf{CP}_k$  (the  $k$ -Covering problem that captures Question (Q1) and (Q2))<sup>3</sup>

**Instance**  $\langle m, n, k \rangle$ :  $U = [1, m]$  and sets  $T_1, T_2, \dots, T_n \subseteq U$  with  $\cup_{i=1}^n T_i = U$ .

**Valid Solutions:** A subset  $U' \subseteq U$  such that  $|U' \cap T_i| \geq n - k$  for each  $i \in [1, n]$ .

**Objective:** *Minimize*  $|U'|$ .

## 2.3 Equivalence of $\mathbf{CP}_k$ and $\mathbf{SC}_{n-k}$

We can establish a 1-1 correspondence between an instance  $\langle m, n, k \rangle$  of  $\mathbf{CP}_k$  and an instance  $\langle n, m, n - k \rangle$  of  $\mathbf{SC}_{n-k}$  by defining  $S_i = \{j \mid i \in T_j\}$  for each  $i \in [1, m]$ . It is easy to verify that  $U'$  is a solution to the instance of  $\mathbf{CP}_k$  if and only if the collection of sets  $S_u$  for each  $u \in U'$  is a solution to the instance of  $\mathbf{SC}_{n-k}$ .

<sup>3</sup> $\mathbf{CP}_{n-1}$  is known as the hitting set problem [7, p. 222].

### 3 Approximation Algorithms for $\mathbf{SC}_k$

An  $\varepsilon$ -approximate solution (or simply an  $\varepsilon$ -approximation) of a minimization problem is defined to be a solution with an objective value no larger than  $\varepsilon$  times the value of the optimum. It is not difficult to see that  $\mathbf{SC}_k$  is NP-complete even when  $k = n - c$  for some constant  $c > 0$ .

#### 3.1 Analysis of Greedy Heuristic for $\mathbf{SC}_k$ for Large $k$

Johnson [8] provides an example in which the greedy heuristic for some instance of  $\mathbf{SC}$  over  $n$  elements has an approximation ratio of at least  $\log_2 n$ . This approach can be generalized to show the following result.

**Lemma 2** *For any fixed  $c > 0$ , the greedy heuristic (as described in Fact 1(b)) has an approximation ratio of at least  $(\frac{1}{2} - o(1)) (\frac{n-c}{8n-2}) \log_2 n = \Omega(\log n)$  for some instance  $\langle n, m, n - c \rangle$  of  $\mathbf{SC}_{n-c}$ .*

**Proof.** We will create an instance of  $\mathbf{SC}_{n-c}$  with  $n = 2\alpha + \gamma \gg c$  where  $\alpha \gg c$  is a sufficiently large positive integer that is also a power of 4 and  $\gamma = 3 + \log_2 \alpha \gg c$  is the least positive integer such that  $2^{\gamma-1} \geq 2\alpha + \gamma - c$ . Notice that by our choice of parameters  $\frac{2^n}{5} < \alpha < \frac{n}{2} \equiv \alpha = \Theta(n)$  and  $1 + \log_2 n < \gamma < 2 + \log_2 n \equiv \gamma = \Theta(\log n)$ . Let  $\mathbf{S} = \{S_1, S_2, \dots, S_{2^{\gamma-1}}\}$  be the collection of all distinct non-empty subsets of  $[2\alpha + 1, 2\alpha + \gamma]$ . Notice that every  $x \in [2\alpha + 1, 2\alpha + \gamma]$  occurs in exactly  $2^{\gamma-1} \geq n - c$  sets in the collection  $\mathbf{S}$ .

A collection of our sets corresponding to an optimal cover of our instance of  $\mathbf{SC}_{n-c}$  will consist of the  $2^{\gamma+1} - 2$  sets  $[1, \alpha] \cup S_i$  and  $[\alpha + 1, 2\alpha] \cup S_i$  for each set  $S_i \in \mathbf{S}$ . Any  $x \in [2\alpha + 1, 2\alpha + \gamma]$  occurs in exactly  $2^\gamma > 2\alpha + \gamma - c = n - c$  of these sets. Also, each  $x \in [1, 2\alpha]$  occurs in exactly  $2^{\gamma-1} \geq 2\alpha + \gamma - c = n - c$  of these sets. Hence, an optimal solution of this instance of  $\mathbf{SC}_{n-c}$  uses at most  $2^{\gamma+1} - 2 < 8n - 2$  sets. Notice that each set in this optimal cover contains at most  $\alpha + \gamma < \frac{n}{2} + 2 + \log_2 n$  elements.

Now we specify another collection of sets which will force the greedy heuristic to use at least  $(n-c) (\frac{1}{2} - o(1)) \log_2 n$  sets. Partition  $[1, \alpha]$  into  $p = 1 + \log_4 \alpha$  disjoint sets  $P_1, P_2, \dots, P_p$  such that  $|P_i| = \lceil \frac{3}{4^i} \alpha \rceil$  for  $i \in [1, p]$ . Observe that  $p > \log_4 n$ . Similarly, partition  $[\alpha + 1, 2\alpha]$  into  $p = 1 + \log_4 \alpha$  disjoint sets  $Q_1, Q_2, \dots, Q_p$  such that  $|Q_i| = \lceil \frac{3}{4^i} \alpha \rceil$  for  $i \in [1, p]$ . Let  $\mathbf{S}' = \{S_1, S_2, \dots, S_{n-c}\} \subseteq \mathbf{S}$ . Now, for each  $P_i \cup Q_i$  and each distinct  $S_j \in \mathbf{S}'$ , create a set  $T_{i,j} = P_i \cup Q_i \cup S_j$ . We claim that greedy will pick the sets  $T_{1,1}, \dots, T_{1,n-c}, T_{2,1}, \dots, T_{2,n-c}, \dots, T_{q,1}, \dots, T_{q,n-c}$  with  $q = (\frac{1}{2} - o(1)) \log_2 n < p$ . This can be shown by induction as follows:

- The greedy must start by picking the sets  $T_{1,1}, \dots, T_{1,n-c}$  in some arbitrary order. Until all these sets have been picked, the unpicked ones have at least  $\frac{3}{4} 2\alpha = \frac{3}{2} \alpha$  elements that have not been covered  $n-c$  times, whereas each set in the optimal cover has at most  $\alpha + \gamma = \alpha + 3 + \log_2 \alpha$  elements and  $\alpha$  is sufficiently large.
- Inductively, suppose that the greedy has picked all sets  $T_{i,j}$  with  $i < q$  when it considers a  $T_{q,r}$  for possible consideration. Obviously  $T_{q,r}$  contains at least  $\frac{3}{4^q} 2\alpha = \frac{6}{4^q} \alpha$  elements that are not yet covered  $n - c$  times. On the other hand, the number of elements that are not yet covered  $n - c$  times in any set from our optimal cover is at most

$$\gamma + \left(1 - \sum_{i=1}^{q-1} \frac{3}{4^i}\right) \alpha = \gamma + \frac{1}{4^{q-1}} \alpha = \gamma + \frac{4}{4^q} \alpha$$

and  $\frac{6}{4^q}\alpha > \gamma + \frac{4}{4^q}\alpha$  provided  $q < \log_4\left(\frac{2\alpha}{\gamma}\right)$ . Since  $\log_4\left(\frac{2\alpha}{\gamma}\right) > \log_4\left(\frac{4n}{5(2+\log_2 n)}\right) > \log_4\left(\frac{4n}{10\log_2 n}\right) = \left(\frac{1}{2} - o(1)\right)\log_2 n$ , the inequality  $\frac{6}{4^q}\alpha > \gamma + \frac{4}{4^q}\alpha$  holds for  $q \in \left[1, \left(\frac{1}{2} - o(1)\right)\log_2 n\right]$ .  $\square$

### 3.2 Randomized Approximation Algorithm for $\mathbf{SC}_k$

As stated before, an instance  $\langle n, m, k \rangle$  of  $\mathbf{SC}_k$  can be  $(1 + \ln \alpha)$ -approximated in  $O(mnk)$  time for any  $k$  where  $\alpha = \max_{S \in \mathcal{S}}\{|S|\}$ . In this section, we provide a randomized algorithm with an expected performance ratio better than  $(1 + \ln \alpha)$  for larger  $k$ . Let  $\mathcal{S} = \{S_1, S_2, \dots, S_m\}$ .

Our algorithm presented below as well as our subsequent discussions and proofs are formulated with the help of the following vector notations:

- All our vectors have  $m$  coordinates with the  $i^{\text{th}}$  coordinate indexed with the  $i^{\text{th}}$  set  $S_i$  of  $\mathcal{S}$ .
- if  $V \subset \mathcal{S}$ , then  $v \in \{0, 1\}^m$  is the characteristic vector of  $V$ , *i.e.*,  $v_{S_i} = \begin{cases} 1 & \text{if } S_i \in V \\ 0 & \text{if } S_i \notin V \end{cases}$
- $\mathbf{1}$  is the vector of all 1's, *i.e.*  $\mathbf{1} = \mathbf{s}$ ;
- $s^i = \{A \in \mathcal{S} : i \in A\}$  denotes the sets in  $\mathcal{S}$  that contains a specific element  $i$ .

Consider the standard integer programming (IP) formulation of an instance  $\langle n, m, k \rangle$  of  $\mathbf{SC}_k$  [18]:

$$\text{minimize } \mathbf{1}x \text{ subject to } \begin{array}{ll} s^i x \geq k & \text{for each } i \in \mathbf{U} \\ x_A \in \{0, 1\} & \text{for each } A \in \mathcal{S} \end{array}$$

A linear programming (LP) relaxation of the above formulation is obtained by replacing each constraint  $x_A \in \{0, 1\}$  by  $0 \leq x_A \leq 1$ . The following randomized approximation algorithm for  $\mathbf{SC}_k$  can then be designed:

1. Select an appropriate positive constant  $\beta > 1$  in the following manner:
$$\beta = \begin{cases} \ln \alpha & \text{if } k = 1 \\ \ln(\alpha/(k-1)) & \text{if } \alpha/(k-1) \geq e^2 \text{ and } k > 1 \\ 2 & \text{if } \frac{1}{4} < \alpha/(k-1) < e^2 \text{ and } k > 1 \\ 1 + \sqrt{\frac{\alpha}{k}} & \text{otherwise} \end{cases}$$
2. Find a solution  $x$  to the LP relaxation via any polynomial-time algorithm for solving linear programs (e.g. [9]).
3. **(deterministic rounding)** Form a family of sets  $\mathcal{C}^0 = \{A \in \mathcal{S} : \beta x_A \geq 1\}$ .
4. **(randomized rounding)** Form a family of sets  $\mathcal{C}^1 \subset \mathcal{S} - \mathcal{C}^0$  by independent random choices such that  $\Pr[A \in \mathcal{C}^1] = \beta x_A$ .
5. **(greedy selection)** Form a family of sets  $\mathcal{C}^2$  as:
  - while  $s^i(c^0 + c^1 + c^2) < k$  for some  $i \in \mathbf{U}$ , insert to  $\mathcal{C}^2$  any  $A \in S^i - \mathcal{C}^0 - \mathcal{C}^1 - \mathcal{C}^2$ .
6. Return  $\mathcal{C} = \mathcal{C}^0 \cup \mathcal{C}^1 \cup \mathcal{C}^2$  as our solution.

Let  $r(\alpha, k)$  denote the performance ratio of the above algorithm.

**Theorem 3**<sup>4</sup>

$$\mathbf{E}[r(a, k)] \leq \begin{cases} 1 + \ln a, & \text{if } k = 1 \\ (1 + e^{-(k-1)/5}) \ln(a/(k-1)), & \text{if } a/(k-1) \geq e^2 \approx 7.39 \text{ and } k > 1 \\ \min\{2 + 2 \cdot e^{-(k-1)/5}, 2 + (e^{-2} + e^{-9/8}) \cdot \frac{a}{k}\} \\ \approx \min\{2 + 2 \cdot e^{-(k-1)/5}, 2 + 0.46 \cdot \frac{a}{k}\} & \text{if } \frac{1}{4} < a/(k-1) < e^2 \text{ and } k > 1 \\ 1 + 2\sqrt{\frac{a}{k}} & \text{if } a/(k-1) \leq \frac{1}{4} \text{ and } k > 1 \end{cases}$$

Let  $\text{OPT}$  denote the minimum number of sets used by an optimal solution. Obviously,  $\text{OPT} \geq \mathbf{1}x$  and  $\text{OPT} \geq \frac{n^k}{a}$ . A proof of Theorem 3 follows by showing the following upper bounds on  $\mathbf{E}[r(a, k)]$  and taking the best of these bounds for each value of  $a/(k-1)$ :

$$\begin{aligned} & 1 + \ln a, & \text{if } k = 1 \\ & (1 + e^{-(k-1)/5}) \ln(a/(k-1)), & \text{if } a/(k-1) \geq e^2 \text{ and } k > 1 \\ & 2 + 2 \cdot e^{-(k-1)/5}, & \text{if } a/(k-1) < e^2 \text{ and } k > 1 \\ & 2 + (e^{-2} + e^{-9/8}) \cdot \frac{a}{k}, & \text{if } a/(k-1) < e^2 \text{ and } k > 1 \\ & 1 + 2\sqrt{\frac{a}{k}}, & \text{if } a/k \leq \frac{1}{2} \end{aligned}$$

**3.2.1 Proof of  $\mathbf{E}[r(a, k)] \leq 1 + \ln a$  if  $k = 1$ ,**

**$\mathbf{E}[r(a, k)] \leq (1 + e^{-(k-1)/5}) \ln(a/(k-1))$  if  $a/(k-1) \geq e^2$  and  $k > 1$ , and**

**$\mathbf{E}[r(a, k)] \leq 2 + 2 \cdot e^{-(k-1)/5}$  if  $a/(k-1) < e^2$  and  $k > 1$**

For our analysis, we use the following notations:

$$x_A^0 = \begin{cases} x_A & \text{if } \beta x_A \geq 1 \\ 0 & \text{otherwise} \end{cases} \quad x_A^1 = \begin{cases} 0 & \text{if } \beta x_A \geq 1 \\ x_A & \text{otherwise} \end{cases}$$

Note that  $c_A^0 = \lceil x_A^0 \rceil \leq \beta x_A^0$ . Thus  $\mathbf{1}x^0 \leq \mathbf{1}c^0 \leq \beta \mathbf{1}x^0$ . Define  $\text{bonus} = \beta \mathbf{1}x^0 - \mathbf{1}c^0$ . It is easy to see that  $\mathbf{E}[\mathbf{1}(c^0 + c^1)] = \beta \mathbf{1}x - \text{bonus}$ .

The contribution of set  $A$  to  $\text{bonus}$  is  $\beta x_A^0 - c_A^0$ . This contribution to  $\text{bonus}$  can be distributed equally to the elements in  $A$ . Since  $|A| \leq a$ , an element  $i \in [1, n]$  receives a total of *at least*  $b^i/a$  of  $\text{bonus}$ , where  $b^i = s^i(\beta x^0 - c^0)$ . The random process that forms set  $\mathcal{C}^1$  has the following goal from the point of view of element  $i$ : pick at least  $g^i$  sets that contain  $i$ , where  $g^i = k - s^i c^0$ . These sets are obtained as successes in Poisson trials whose probabilities of success add to at least  $p^i = \beta(k - s^i x^0)$ . Let  $y^i$  be random function denoting the number that element  $i$  contributes to the size of  $\mathcal{C}^2$ ; thus, if in the random trials in Step 4 we found  $h$  sets from  $S^i$  then  $y^i = \max\{0, k - h\}$ . Thus,  $\mathbf{E}[r(a, k)] = \mathbf{E}[\mathbf{1}(c^0 + c^1 + c^2)] \leq \beta \mathbf{1}x + \sum_{i=1}^n \mathbf{E}[y^i - \frac{b^i}{a}]$ . Let  $q^i = \frac{\beta}{\beta-1} s^i(c^0 - x^0)$ . We can parameterize the random process that forms the set  $\mathcal{C}^2$  from the point of view of element  $i$  as follows:

- $g^i$  is the *goal* for the number of sets to be picked;
- $p^i = \beta(k - s^i x^0) = \beta g^i + (\beta - 1)q^i$  is the sum of probabilities with which sets are picked;
- $b^i/a$  is the *bonus* of  $i$ , where  $b^i = s^i(\beta x^0 - c^0) \geq (\beta - 1)(k - g^i - q^i)$ ;

<sup>4</sup>The case of  $k = 1$  was known before and included for the sake of completeness only.

- $q^i \geq 0$ ,  $g^i \geq 0$  and  $g^i + q^i \leq k$ ;
- $y^i$  measures how much the goal is *missed*;
- to bound  $\mathbf{E}[r(a, k)]$  we need to bound  $\mathbf{E}[y^i - \frac{b^i}{a}]$ .

### 3.2.1.1 g-shortage Functions

In this section we prove some inequalities needed to estimate  $\mathbf{E}[y^i - \frac{b^i}{a}]$  tightly. Assume that we have a random function  $X$  that is a sum of  $N$  independent 0-1 random variables  $X_i$ . Let  $\mathbf{E}[X] = \sum_i \mathbf{Pr}[X_i = 1] = \mu$  and  $g < \mu$  be a positive integer. We define *g-shortage function* as  $Y_g^\mu = \max\{g - X, 0\}$ . Our goal is to estimate  $\mathbf{E}[Y_g^\mu]$ .

**Lemma 4**  $\mathbf{E}[Y_g^\mu] < e^{-\mu} \sum_{i=0}^{g-1} \frac{g-i}{i!} \mu^i$ .

**Proof.** Suppose that for some positive numbers  $p, q$  and some  $X_i$  we have  $\mathbf{Pr}[X_i = 1] = p + q$ . Consider replacing  $X_i$  with two independent random functions  $X_{i,0}$  and  $X_{i,1}$  such that  $\mathbf{Pr}[X_{i,0} = 1] = p$  and  $\mathbf{Pr}[X_{i,1} = 1] = q$ . We can show that after this replacement  $\mathbf{E}[Y_g^\mu]$  increases as follows. In terms of our random functions *before the replacement* we define  $r_j = \mathbf{Pr}[X - X_i = g - j]$ . Let  $X'$  be the sum of our random functions *after the replacement* and  $Y'_g$  be defined in terms of  $X'$ . Let  $a = \sum_{j=1}^{g-1} jr_j$ ,  $b = \sum_{j=2}^{g-1} (j-1)r_j$ , and  $c = \sum_{j=3}^{g-1} (j-2)r_j$ . Then,

$$\begin{aligned} \mathbf{E}[Y'_g] &= (1-p)(1-q)a + p(1-q)b + (1-p)qb + pqc \\ &= (1-p-q+pq)a + (p+q-2pq)b + pqc \\ \mathbf{E}[Y_g] &= (1-p-q)a + (p+q)b \\ \mathbf{E}[Y'_g] - \mathbf{E}[Y_g] &= (a-2b+c)pq = r_1pq \geq 0 \end{aligned}$$

Therefore, we increase  $\mathbf{E}[Y_g^\mu]$  if we replace the original independent random function by  $N$  *Bernoulli trials* with probability of success  $\mu/N$ , and take the limit for  $N \rightarrow \infty$ . If  $r_j = \mathbf{Pr}[X_g = j]$  then it now follows that

$$\lim_{N \rightarrow \infty} r_j = \lim_{N \rightarrow \infty} \frac{N!}{(N-j)!j!} \left(1 - \frac{\mu}{N}\right)^{N-j} \left(\frac{\mu}{N}\right)^j = \frac{\mu^j}{e^{\mu j}}$$

where the last equality follows from standard estimates in probability theory.  $\square$

From now on we will assume the worst-case distribution of  $Y_g^\mu$ , so we will assume that the above inequality in Lemma 4 is actually an equality (as it becomes so in the limit), *i.e.*, we assume  $\mathbf{E}[Y_g^\mu] = e^{-\mu} \sum_{i=0}^{g-1} \frac{g-i}{i!} \mu^i$ . For a fixed  $\beta$ , we will need to estimate the growth of  $\mathbf{E}[Y_g^{g\beta}]$  as a function of  $g$ . Let  $\rho_g(\beta) = e^{g\beta} \mathbf{E}[Y_g^{g\beta}]$ .

**Lemma 5**  $\rho_g(1) = \sum_{i=0}^{g-1} \frac{g-i}{i!} g^i = \frac{g^g}{(g-1)!}$

**Proof.**

$$\begin{aligned} \rho_g(1) &= \sum_{i=0}^{g-1} \frac{g^i(g-i)}{i!} = \sum_{i=0}^{g-1} \frac{g^{i+1}}{i!} - \sum_{i=0}^{g-1} \frac{g^i}{i!} \\ &= \sum_{i=0}^{g-1} \frac{g^{i+1}}{i!} - \sum_{i=1}^{g-1} \frac{g^i}{(i-1)!} \end{aligned}$$

$$\begin{aligned}
&= \sum_{i=0}^{g-1} \frac{g^{i+1}}{i!} - \sum_{i=0}^{g-2} \frac{g^{i+1}}{i!} \\
&= \sum_{i=0}^{g-1} \frac{g^{i+1}}{i!} - \left( \sum_{i=0}^{g-1} \frac{g^{i+1}}{i!} - \frac{g^g}{(g-1)!} \right) \\
&= \frac{g^g}{(g-1)!}
\end{aligned}$$

□

**Lemma 6** For  $\beta > 1$ ,  $\frac{\rho_{g+1}(\beta)}{\beta \rho_g(\beta)}$  is a decreasing function of  $\beta$ .

**Proof.** By definition,  $\rho_g(\beta) = \sum_{i=0}^{g-1} \frac{\alpha_i}{i!} \beta^i$  where  $\alpha_i = g^i(g-i)$ . Let  $f(\beta) = \rho_{g+1}(\beta)$  and  $t(\beta) = \beta \rho_g(\beta)$ . We need to show that, for a given fixed  $g$ ,  $f(\beta)/t(\beta)$  is a decreasing function of  $\beta$ . The derivative of  $f(\beta)/t(\beta)$  is  $\frac{f'(\beta)t(\beta) - t'(\beta)f(\beta)}{(t(\beta))^2}$ . We claim that the numerator  $f'(\beta)t(\beta) - t'(\beta)f(\beta)$  is a polynomial whose coefficients are *always negative*, which then proves that the derivative is negative (for all  $\beta > 0$ ). To prove this, it suffices to show that if  $p(\beta) = f'(\beta)t(\beta) - t'(\beta)f(\beta)$  then  $p^{(k)}(0) < 0$  for all  $k$ .

Note that  $t^{(k)}(0) = k \rho_g^{(k-1)}(0)$  for all  $k$ , hence

$$f^{(k)}(0) = \begin{cases} (g+1)^k(g+1-k), & \text{if } 0 \leq k \leq g \\ 0, & \text{if } k > g \end{cases} \quad t^{(k)}(0) = \begin{cases} kg^{k-1}(g+1-k), & \text{if } 0 \leq k \leq g \\ 0, & \text{if } k > g \end{cases}$$

On the other hand,

$$\begin{aligned}
p' &= f't - t'f, \\
p'' &= f''t - t''f, \\
p''' &= f'''t - t'''f + f''t' - t''f'
\end{aligned}$$

etc, so it will be enough to prove that

$$a(k, h) = f^{(k)}t^{(h)}(0) - t^{(k)}f^{(h)}(0) < 0$$

whenever  $g \geq k > h$ . (By induction, if we have that a derivative of  $p$  is a sum of terms of the form  $t^{(k)}f^{(h)} - f^{(k)}t^{(h)}$ , then taking another derivative of any such term, we get 4 terms, then rearrange to see that the same is true for one more derivative of  $p$ ).

Let  $\kappa = h + d$  with  $d > 0$ . Then  $a(h + d, h) = (-g^d h - g^d d + (g+1)^d h)(g+1)^h g^{h-1}$ , so we need to show that  $(g+1)^d h < g^d(h+d)$  when  $g \geq h+d$  and  $d > 0$ . Let  $q(h) = -g^d(h+d) + (g+1)^d h$  for fixed  $g$  and  $d$ . We need to show that  $q(h) < 0$  for  $h \leq g-d$ . Since  $q'(h) = (g+1)^d - g^d > 0$  it is enough to check at the maximum value  $h = g-d$ . Thus, we need to show that  $(g+1)^d(g-d) < g^{d+1}$  always holds for any  $1 \leq d \leq g$ . For  $d = g$  this is clear, so we assume from now on that  $d < g$ . Taking logarithms of both sides, we must show that  $r(d) = d \ln(g+1) + \ln(g-d) - (d+1) \ln g < 0$ . We claim that  $r'(d) = \ln\left(\frac{g+1}{g}\right) - \frac{1}{g-d} < 0$ . Once that this is shown, it will follow from  $r(d) < r(1) = \ln(g+1) + \ln(g-1) - 2 \ln g$  and concavity of  $\ln$  (which says that  $r(1) < 0$ ) that  $r(d) < 0$  as wanted. To show  $r'(d) < 0$ , we note that  $-\frac{1}{g-d} < -\frac{1}{g-1}$  (because  $d \geq 1$ ), so all we need to show is that  $\ln\left(\frac{g+1}{g}\right) < \frac{1}{g-1}$ , or, equivalently,  $1 + \frac{1}{g} < e^{\frac{1}{g-1}}$ . But, obviously,  $1 + \frac{1}{g} < 1 + \frac{1}{g-1} < e^{\frac{1}{g-1}}$ . □

The next lemma characterizes the growth of  $\mathbf{E}[Y_g^{gX}]$  as a function of  $g$ .

**Lemma 7** If  $g > 1$  and  $\beta > 1$  then  $\frac{\mathbf{E}[Y_g^{g\beta}]}{\mathbf{E}[Y_{g-1}^{(g-1)\beta}]} \leq e^{-\beta} \left(\frac{g}{g-1}\right)^g$

**Proof.** Using Lemma 5 and Lemma 6, we get  $\frac{\mathbf{E}[Y_g^{g\beta}]}{\mathbf{E}[Y_{g-1}^{(g-1)\beta}]} = e^{-\beta} \beta \frac{\rho_g(\beta)}{\beta \rho_{g-1}(\beta)} \leq e^{-\beta} \beta \frac{\rho_g(1)}{\rho_{g-1}(1)} = e^{-\beta} \beta \left(\frac{g}{g-1}\right)^g$ .  $\square$

The last lemma characterizes the impact of “extra probabilities” on the expected value.

**Lemma 8**  $\frac{\mathbf{E}[Y_g^{g\beta+q}]}{\mathbf{E}[Y_g^{g\beta}]} < e^{-q(1-1/\beta)}$

**Proof.** The ratio is  $e^{-q}$  times the ratio of two polynomials. The terms of the upper polynomial are larger than the terms of the lower by a factor at most  $\left(\frac{g\beta+q}{g\beta}\right)^{g-1} = \left(1 + \frac{q}{g\beta}\right)^{g-1} < e^{q/\beta}$ .  $\square$

### 3.2.1.2 Putting All the Pieces Together

In this section we put all the pieces together from the previous two subsections to prove our claim on  $\mathbf{E}[r(\mathbf{a}, k)]$ . We assume that  $\beta \geq 2$  if  $k > 1$ . Because we perform analysis from the point of view of a fixed element  $i$ , we will skip  $i$  as a superscript as appropriate. As we observed in Section 3.2.1, we need to estimate  $\mathbf{E}[y - \frac{b}{a}]$  and  $b \geq (\beta - 1)(k - g - q)$ . We will also use the notations  $p$  and  $q$  as defined there.

Obviously if  $g = 0$  then  $y = 0$ . First, we consider the case of  $k = 1$  separately. Then,  $g \leq 1$  and hence  $\mathbf{E}[y - \frac{b}{a}] \leq e^{-\beta} = \frac{1}{a}$ . Thus, when  $k = 1$ ,

$$\begin{aligned} \mathbf{E}[r(\mathbf{a}, k)] &= \mathbf{E}[1(c^0 + c^1 + c^2)] \\ &\leq \beta \mathbf{1}x + \sum_{i=1}^n \mathbf{E}[y^i - \frac{b^i}{a}] \\ &\leq \beta \mathbf{1}x + \frac{n}{a} \\ &\leq \beta \text{OPT} + \text{OPT} = (1 + \ln a) \text{OPT} \end{aligned}$$

Otherwise, for the rest of this proof, assume that  $k > 1$ . We first consider the “base” case of  $g = 1$  and  $q = 0$ . Since  $q = 0$ ,  $c^0 = x^0$ . Thus,  $b = s^i(\beta c^0 - c^0) = (\beta - 1)s^i c^0 = (\beta - 1)(k - 1)$ . Next, we compute  $\mathbf{E}[y]$ . Since  $p = \beta g = \beta$ ,  $\mathbf{E}[y] = \mathbf{E}[Y_1^\beta] = e^{-\beta}$ .

We postulate that

$$\begin{aligned} \mathbf{E}[y - \frac{b}{a}] \leq 0 &\equiv e^{-\beta} \leq \frac{(\beta - 1)(k - 1)}{a} \\ &\equiv \frac{e^{-\beta}}{\beta - 1} \leq \frac{k - 1}{a} \\ &\equiv e^\beta (\beta - 1) \geq \frac{a}{k - 1} \\ &\equiv \beta + \ln(\beta - 1) \geq \ln \frac{a}{k - 1} \end{aligned} \tag{4}$$

Now we observe the following:

- If  $a/(k - 1) \geq e^2$ , then  $\beta + \ln(\beta - 1) \geq 2$ , or equivalently,  $\beta \geq 2$  as is the assumption in this section. Moreover,  $\beta = \ln(a/(k - 1))$  obviously satisfies inequality 4.

- If  $\alpha/(k-1) < e^2$ , then obviously  $\beta \geq 2$  by our choice of  $\beta = 2$ . Moreover,  $\beta = 2$  obviously also satisfies inequality 4 as well.

Thus, for the base case,  $\mathbf{E}[\mathbf{1}(c^0 + c^1 + c^2)] \leq \beta \mathbf{1}x \leq \ln(\alpha/(k-1))\text{OPT}$ .

Now we consider the “non-base” case when either  $g > 1$  or  $q > 0$ . Compared to the base case, in a non-base case we have bonus  $\frac{b}{a}$  decreased by at least  $(\beta-1)(g+q-1)/\alpha$ . Also,  $\mathbf{E}[y] = \mathbf{E}[Y_g^\beta] = \mathbf{E}[Y_g^{\beta g + (\beta-1)q}]$ . We need to calculate how this compares with the base value of  $\mathbf{E}[Y_1^\beta]$  using Lemma 7 and Lemma 8.

**Lemma 9**  $\frac{\mathbf{E}[Y_g^{\beta g + (\beta-1)q}]}{\mathbf{E}[Y_1^\beta]} \leq e^{-(g+q-1)/5}$ .

**Proof.** Firstly, if  $q > 0$ , then

$$\begin{aligned} \frac{\mathbf{E}[Y_g^{\beta g + (\beta-1)q}]}{\mathbf{E}[Y_g^{\beta g}]} &< e^{-(\beta-1)q\left(1-\frac{1}{\beta}\right)} \quad (\text{by Lemma 8}) \\ &\leq e^{-(\beta-1)q\left(1-\frac{1}{2}\right)} \quad (\text{since } \beta \geq 2) \\ &\leq e^{-q/2} \quad (\text{since } \beta \geq 2) \\ &< e^{-q/5} \end{aligned}$$

Now we need to bound  $\mathbf{E}[Y_g^{\beta g}]/\mathbf{E}[Y_1^\beta]$ .

Obviously,  $\frac{\mathbf{E}[Y_g^{\beta g}]}{\mathbf{E}[Y_1^\beta]} = \prod_{i=2}^g \frac{\mathbf{E}[Y_i^{\beta i}]}{\mathbf{E}[Y_{i-1}^{\beta(i-1)}]}$ . We now observe the following:

- If  $i = 2$ , then  $\frac{\mathbf{E}[Y_2^{\beta 2}]}{\mathbf{E}[Y_1^{\beta 1}]} = e^{-\beta \frac{\rho_2(\beta)}{\rho_1(\beta)}} = e^{-\beta(2+2\beta)}$ . Since  $e^{-\beta(2+2\beta)}$  is a decreasing function of  $\beta$  and  $2+\beta > 2$ , it follows that  $e^{-\beta(2+2\beta)} < 4e^{-2} < e^{-1/5}$ .
- Similarly, if  $i = 3$ , then  $\frac{\mathbf{E}[Y_3^{\beta 3}]}{\mathbf{E}[Y_2^{\beta 2}]} = e^{-2\beta \frac{\rho_3(\beta)}{\rho_2(\beta)}} = e^{-2\beta(3+6\beta+\frac{9}{2}\beta^2)} < 33e^{-4} < e^{-2/5}$ .
- Similarly, if  $i = 4$ , then  $\frac{\mathbf{E}[Y_4^{\beta 4}]}{\mathbf{E}[Y_3^{\beta 3}]} = e^{-3\beta \frac{\rho_4(\beta)}{\rho_3(\beta)}} = e^{-3\beta(4+12\beta+16\beta^2+\frac{32}{3}\beta^3)} < \frac{532}{3}e^{-6} < e^{-3/5}$ .
- Similarly, if  $i = 5$ , then  $\frac{\mathbf{E}[Y_5^{\beta 5}]}{\mathbf{E}[Y_4^{\beta 4}]} = e^{-4\beta \frac{\rho_5(\beta)}{\rho_4(\beta)}} = e^{-4\beta(5+20\beta+\frac{75}{2}\beta^2+\frac{125}{3}\beta^3+\frac{625}{24}\beta^4)} < 945e^{-8} < e^{-4/5}$ .
- Finally, suppose that  $i \geq 6$ . Then,

$$\begin{aligned} \frac{\mathbf{E}[Y_i^{\beta i}]}{\mathbf{E}[Y_{i-1}^{\beta(i-1)}]} &\leq e^{-\beta \beta \left(\frac{i}{i-1}\right)^i} \quad (\text{by Lemma 7}) \\ &< e^{-\beta \beta \left(\frac{6}{5}\right)^6} \quad (\text{since } \left(\frac{i}{i-1}\right)^i \text{ is a decreasing function of } i) \\ &\leq e^{-22 \left(\frac{6}{5}\right)^6} \quad (\text{since } e^{-\beta \beta} \text{ is a decreasing function of } \beta) \\ &< e^{-1/5} \end{aligned}$$

- Thus,  $\frac{\mathbf{E}[Y_g^{\beta g}]}{\mathbf{E}[Y_1^\beta]} \leq e^{-(g-1)/5}$ .
- Thus,  $\frac{\mathbf{E}[Y_g^{\beta g + (\beta-1)q}]}{\mathbf{E}[Y_1^\beta]} \leq e^{-(g+q-1)/5}$ .

□

Summarizing, when bonus is decreased by at most  $(\beta-1)(g+q-1)/a = (\beta-1)t/a$ , we decrease the estimate of  $\mathbf{E}[y]$  by multiplying it with at least  $e^{-t/5}$ . As a function of  $t = g + q - 1$  we have

$$\mathbf{E}[y] - b/a \leq e^{-\beta-t/5} - \frac{\beta-1}{a}(k-1-t) \leq \frac{(\beta-1)(k-1)}{a} \left( e^{-t/5} - 1 + \frac{t}{k-1} \right)$$

This is a convex function of  $t$ , so its maximal value must occur at one of the ends of its range. When  $t = 0$  we have 0, and when  $t = k - 1$  we have  $\frac{(\beta-1)(k-1)}{a}e^{-(k-1)/5}$ . As a result, our expected performance ratio for  $k > 1$  is given by

$$\begin{aligned} \mathbf{E}[r(a, k)] &\leq \beta \mathbf{1}_x + \sum_{i=1}^n \mathbf{E}[y^i - \frac{b^i}{a}] \\ &\leq \beta \text{OPT} + \frac{\beta n k}{a} e^{-(k-1)/5} \\ &\leq \beta(1 + e^{-(k-1)/5}) \text{OPT} \\ &\leq \begin{cases} (1 + e^{-(k-1)/5}) \ln(a/(k-1)) \text{OPT} & \text{if } a/(k-1) \geq e^2 \\ 2 \cdot (1 + e^{-(k-1)/5}) \text{OPT} & \text{if } a/(k-1) < e^2 \end{cases} \end{aligned}$$

### 3.2.2 Proof of $\mathbf{E}[r(a, k)] \leq 2 + (e^{-2} + e^{-9/8}) \cdot \frac{a}{k}$ if $a/(k-1) < e^2$

Each set in  $A \in \mathcal{C}_0 \cup \mathcal{C}_1$  is selected with probability  $\min\{\beta x_A, 1\}$ . Thus  $\mathbf{E}[|\mathcal{C}_0 \cup \mathcal{C}_1|] \leq \beta \mathbf{1}_x \leq \beta \cdot \text{OPT}$ . Next we estimate an upper bound on  $\mathbf{E}[|\mathcal{C}_2|]$ . For each element  $i \in [1, n]$  let the random variable  $v_i$  be  $\max\{k - d_i, 0\}$  where  $d_i$  is the number of sets in  $\mathcal{C}_0 \cup \mathcal{C}_1$  that contain  $i$ . Clearly,  $|\mathcal{C}_2| \leq \sum_{i=1}^n v_i$ . Thus it suffices to estimate  $\mathbf{E}[v_i]$ . Because our estimate will not depend on  $i$ , we will drop this index  $i$  from  $v_i$  for notational simplifications. Assume that  $i \in [1, n]$  is contained in  $k - f$  sets from  $\mathcal{C}_0$  for some  $f < k$ . Then,  $1 \leq v \leq f$  and

$$\mathbf{E}[v] \leq \sum_{j=1}^{\infty} \Pr[v \geq j] = \sum_{j=1}^f \Pr[v \geq j] = \sum_{j=0}^{f-1} \Pr[v \geq f - j] \quad (5)$$

Let the random variable  $y$  denote the number of sets in  $\mathcal{C}_1$  that contain  $i$ . Considering the constraint  $s^i x \geq k$  and the fact that we select a set  $A \in \mathcal{S} \setminus \mathcal{C}_0$  with probability  $\beta x_A$ , it follows that  $\mathbf{E}[y] \geq \beta f$ . Now,

$$\Pr[v \geq f - j] = \Pr[y < (1 - \delta_j)\beta f] \quad (6)$$

where

$$(1 - \delta_j)\beta f = j \quad \equiv \quad \beta f \delta_j = \beta f - j \quad \equiv \quad \delta_j = \frac{\beta f - j}{\beta f} \quad (7)$$

By using standard Chernoff's bound [1, 3, 12], we have

$$\Pr[y < (1 - \delta_j)\beta f] \leq \Pr[y < (1 - \delta_j)\mathbf{E}[y]] < e^{-\mathbf{E}[y]\delta_j^2/2} \leq e^{-\beta f \delta_j^2/2} \quad (8)$$

where

$$\frac{\beta f \delta_j^2}{2} = \frac{(\beta f - j)^2}{2\beta f} = \frac{f}{2}\beta - j + \frac{j^2}{2\beta f} = \zeta(\beta, f, j) \quad (9)$$

Combining Equations (5), (6), (8) and (9) we get  $\mathbf{E}[v] < \sum_{j=0}^{f-1} e^{-\zeta(\beta, f, j)}$ .

**Lemma 10** Let  $X(\beta, f) = \sum_{j=0}^{f-1} e^{-\zeta(\beta, f, j)}$ . Then  $X(2, f)$  is maximized for  $f = 2$  and this maximum value is  $e^{-2} + e^{-9/8}$ .

**Proof.**  $X(2, 1) = e^{-1} < e^{-2} + e^{-9/8}$ . The following series of arguments shows that  $X(2, f) \leq X(2, 2)$  for all  $f > 2$ :

- $\zeta(2, f, j) = f - j + \frac{j^2}{4f}$ . Also,  $\zeta(2, f, f-1) = 1 + \frac{(f-1)^2}{4f} = \frac{f+2+1/f}{4}$  is an increasing function of  $f$ .
- $\zeta(2, p+1, j) - \zeta(2, p, j) = 1 + \frac{j^2}{4} \left( \frac{1}{p+1} - \frac{1}{p} \right) = 1 - \frac{j^2}{4p(p+1)} > \frac{3}{4}$  for all  $p \geq 1$ .
- Clearly,

$$\begin{aligned} X(2, f) &= \left( \sum_{j=0}^{f-2} e^{-\zeta(2, f-1, j)} \cdot e^{\zeta(2, f-1, j) - \zeta(2, f, j)} \right) + e^{-\zeta(2, f, f-1)} \\ &< e^{-3/4} X(2, f-1) + e^{-\zeta(2, f, f-1)} < \frac{1}{2} X(2, f-1) + e^{-\zeta(2, f, f-1)} \end{aligned}$$

Hence

$$\begin{aligned} X(2, f-1) \leq X(2, 2) \wedge e^{-\zeta(2, f, f-1)} < \frac{1}{2} X(2, 2) &\implies X(2, f) < X(2, 2) \\ \equiv \\ X(2, f-1) \leq X(2, 2) \wedge \zeta(2, f, f-1) > -\ln X(2, 2) + \ln 2 &\implies X(2, f) < X(2, 2) \end{aligned}$$

- $X(2, 3) < 0.44 < X(2, 2)$ .
- $\zeta(2, 4, 3) > 1.56 > 0.78 + 0.694 > -\ln X(2, 2) + \ln 2$ . Moreover, since  $\zeta(2, f, f-1)$  increases with  $f$ , this implies  $\zeta(2, f, f-1) > -\ln X(2, 2) + \ln 2$  for all  $f \geq 4$ . Thus,  $X(2, f) < X(2, 2)$  for all  $f \geq 4$ .  $\square$

Now we are able to complete the proof on our claimed expected performance ratios as follows. If  $\alpha/(k-1) \leq e^2$  then with  $\beta = 2$  we get  $\mathbf{E}[|\mathcal{C}_0 \cup \mathcal{C}_1|] \leq 2 \cdot \text{OPT}$  and  $\mathbf{E}[|\mathcal{C}_2|] \leq (e^{-2} + e^{-9/8}) \cdot \alpha \leq (e^{-2} + e^{-9/8}) \cdot \frac{\alpha}{k} \cdot \text{OPT}$  by Lemma 10.

### 3.2.3 Proof of $\mathbf{E}[r(\alpha, k)] \leq 1 + 2\sqrt{\frac{\alpha}{k}}$ if $\alpha/k \leq \frac{1}{2}$

For notational simplification, let  $\alpha = \sqrt{\frac{k}{a}} \geq 2$ . Thus, in our notation,  $\beta = 1 + \alpha^{-1}$  and we need to show that  $\mathbf{E}[r(\alpha, k)] \leq 1 + 2\alpha^{-1} + \alpha^{-2}$ . As we observed immediately after the statement of Theorem 3,  $\text{OPT} \geq \mathbf{1x}$  (where  $\mathbf{x}$  was the solution vector to the LP relaxation) and  $\text{OPT} \geq \frac{nk}{a} = n\alpha^2$ . We will also reuse, if necessary, the notations introduced in Section 3.2.1.

We first focus our attention on a single element, say  $i$ . For notational convenience, we will drop  $i$  from the superscript when possible. Let  $\mathcal{C}_0^i$  and  $\mathcal{C}_1^i$  be the sets in  $\mathcal{C}_0$  and  $\mathcal{C}_1$ , respectively, that contained  $i$ . We will relate the following quantities:

- $p = p^i = \beta(k - s^i x^0)$  is the sum of probabilities used in Step 4 for the sets that contain  $i$ ;
- $y = y^i = |\mathcal{C}_2^i|$  is the shortage of sets that contain  $i$  after Step 4.

Suppose that  $y > 0$  and that  $|\mathcal{C}_0^i| = k - f$  for some  $0 < f \leq k$ . Suppose that element  $i$  is contained in  $k - f + \rho$  sets, say the sets  $S_1, S_2, \dots, S_{k-f+\rho}$ , for some  $\rho > 0$ , out of which  $k - f$  sets, say the sets  $S_1, S_2, \dots, S_{k-f}$ , were selected to be in  $\mathcal{C}_0^i$ . From the inequality

$$x_{S_1} + x_{S_2} + \dots + x_{S_{k-f+\rho}} \geq k$$

and the fact that  $x_j \leq 1$  for all  $j$ , it follows that  $\rho \geq x_{S_{k-f+1}} + x_{S_{k-f+2}} + \dots + x_{S_{k-f+\rho}} \geq f$ . Let  $\rho = f + \mu$  for some  $\mu \geq 0$ . Obviously,  $\mathbf{E}[|\mathcal{C}_1^i|] = p = \beta\rho = (1 + \alpha^{-1})f + (1 + \alpha^{-1})\mu$ . Now,

$$\begin{aligned}
|\mathcal{C}_1^i| &= f - y \\
&= [(1 + \alpha^{-1})f + (1 + \alpha^{-1})\mu] - \alpha^{-1}f - (1 + \alpha^{-1})\mu - y \\
&= \left[ 1 - \left( \frac{\alpha^{-1}f + (1 + \alpha^{-1})\mu}{(1 + \alpha^{-1})f + (1 + \alpha^{-1})\mu} \right) - \frac{y}{(1 + \alpha^{-1})f + (1 + \alpha^{-1})\mu} \right] p \\
&< \left[ 1 - \left( \frac{\alpha^{-1}f}{(1 + \alpha^{-1})f} \right) - \frac{y}{(1 + \alpha^{-1})f + (1 + \alpha^{-1})\mu} \right] p \\
&= \left[ 1 - (1 + \alpha)^{-1} - \frac{y}{p} \right] \mathbf{E}[|\mathcal{C}_1^i|]
\end{aligned}$$

By using standard Chernoff's bound [1, 3, 12], we have

$$\begin{aligned}
\mathbf{E}[|\mathcal{C}_2^i|] &= \sum_{j=1}^{\infty} \mathbf{Pr}[y \geq j] \\
&< \sum_{j=1}^{\infty} \mathbf{Pr}[|\mathcal{C}_1^i| \leq \left(1 - \left((1 + \alpha)^{-1} + \frac{j}{p}\right)\right)] \mathbf{E}[|\mathcal{C}_1^i|] \\
&\leq \sum_{j=1}^{\infty} e^{-\frac{1}{2}((1 + \alpha)^{-1} + j/p)^2 p} \\
&< \sum_{j=1}^{\infty} e^{-\frac{1}{2}(4j(1 + \alpha)^{-1}/p)} \quad \text{since } (x + y)^2 \geq 4xy \text{ for all } x \text{ and } y \\
&= \sum_{j=1}^{\infty} e^{-\frac{1}{2}(4j(1 + \alpha)^{-1})} \\
&= \sum_{j=1}^{\infty} \left[ e^{2(1 + \alpha)^{-1}} \right]^{-j} \\
&= \frac{1}{e^{2(1 + \alpha)^{-1}} - 1} \\
&\leq \frac{1}{2(1 + \alpha)^{-1}} \quad \text{since } e^x > 1 + x \text{ for } x > 0 \\
&= \frac{1 + \alpha}{2}
\end{aligned}$$

Therefore on behalf of  $i$  we will select, on average,  $(1 + \alpha)/2$  sets in Step 5; thus the total average number of elements selected in Step 5 over all elements is at most  $n\alpha/2$ . Summarizing,

$$\mathbf{E}[|\mathcal{C}^0 + \mathcal{C}^1 + \mathcal{C}^2|] \leq 1x(1 + \alpha^{-1}) + \frac{n\alpha^2}{2\alpha} + \frac{n\alpha^2}{2\alpha^2} \leq \left(1 + \frac{3}{2\alpha} + \frac{1}{\alpha^2}\right) \text{OPT} \leq \left(1 + \frac{2}{\alpha}\right) \text{OPT}$$

where the last inequality follows since  $\alpha \geq 2$ .

## Acknowledgements

We would like to thank Uriel Feige for his constructive criticisms of an earlier draft of the manuscript that led to the improved bounds reported in Section 3.2.3. We also thank the reviewers for their constructive comments.

## References

- [1] N. Alon and J. Spencer, *The Probabilistic Method*, Wiley Interscience, New York, 1992.
- [2] M. Andrec, B.N. Kholodenko, R.M. Levy, and E.D. Sontag. *Inference of signaling and gene regulatory networks by steady-state perturbation experiments: Structure and accuracy*, J. Theoretical Biology, in press.

- [3] H. Chernoff. *A measure of asymptotic efficiency of tests of a hypothesis based on the sum of observations*, Annals of Mathematical Statistics, 23: 493–509, 1952.
- [4] T. Cover. *Geometrical and statistical properties of systems of linear inequalities with applications in pattern recognition*, IEEE Trans. Electronic Computers, EC-14, pp. 326–334, 1965. Reprinted in *Artificial Neural Networks: Concepts and Theory*, IEEE Computer Society Press, Los Alamitos, Calif., 1992, P. Mehra and B. Wah, eds.
- [5] E.J. Crampin, S. Schnell, and P.E. McSharry. *Mathematical and computational techniques to deduce complex biochemical reaction mechanisms*, Progress in Biophysics & Molecular Biology, 86, pp. 77-112, 2004.
- [6] U. Feige. *A threshold for approximating set cover*, Journal of the ACM, Vol. 45, 1998, pp. 634-652.
- [7] M. R. Garey and D. S. Johnson. *Computers and Intractability - A Guide to the Theory of NP-Completeness*, W. H. Freeman & Co., 1979.
- [8] D. S. Johnson. *Approximation Algorithms for Combinatorial Problems*, Journal of Computer and Systems Sciences, Vol. 9, 1974, pp. 256-278.
- [9] N. Karmarkar. *A new polynomial-time algorithm for linear programming*, Combinatorica, 4: 373–395, 1984.
- [10] B. N. Kholodenko, A. Kiyatkin, F. Bruggeman, E.D. Sontag, H. Westerhoff, and J. Hoek. *Untangling the wires: a novel strategy to trace functional interactions in signaling and gene networks*, Proceedings of the National Academy of Sciences USA 99, pp. 12841-12846, 2002.
- [11] B. N. Kholodenko and E.D. Sontag. *Determination of functional network structure from local parameter dependence data*, arXiv physics/0205003, May 2002.
- [12] R. Motwani and P. Raghavan. *Randomized Algorithms*, Cambridge University Press, New York, NY, 1995.
- [13] R. Raz and S. Safra. *A sub-constant error-probability low-degree test and sub-constant error-probability PCP characterization of NP*, proceedings of the 29th Annual ACM Symposium on Theory of Computing, pp. 475-484, 1997.
- [14] L. Schläfli. *Theorie der vielfachen Kontinuität (1852)*, in *Gesammelte Mathematische Abhandlungen*, volume 1, pp. 177–392, Birkhäuser, Basel, 1950.
- [15] E. D. Sontag. *VC dimension of neural networks*, in *Neural Networks and Machine Learning* (C.M. Bishop, ed.), Springer-Verlag, Berlin, pp. 69-95, 1998.
- [16] E.D. Sontag, A. Kiyatkin, and B.N. Kholodenko. *Inferring dynamic architecture of cellular networks using time series of gene expression, protein and metabolite data*, Bioinformatics 20, pp. 1877-1886, 2004.
- [17] J. Stark, R. Callard and M. Hubank. *From the top down: towards a predictive biology of signaling networks*, Trends Biotechnol. 21, pp. 290-293, 2003.
- [18] V. Vazirani. *Approximation Algorithms*, Springer-Verlag, July 2001.