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

Piotr Berman[1], Bhaskar DasGupta[2], and Eduardo Sontag[3]

[1]  Department of Computer Science and Engineering, Pennsylvania State University, University Park, PA 16802. Email: `berman@cse.psu.edu`.
[2]  Department of Computer Science, University of Illinois at Chicago, Chicago, IL 60607-7053. Email: `dasgupta@cs.uic.edu`.
[3]  Department of Mathematics, Rutgers University, New Brunswick, NJ 08903. Email: `sontag@hilbert.rutgers.edu`.

**Abstract.** In this paper we investigate the computational complexities 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 $1 < a < n$ denotes the maximum number of elements in any given set in our set multicover problem. Then, we show that a nontrivial 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 $\mathbf{e}^2 \approx 7.39$, and for smaller values of $a/k$ it decreases towards 2 exponentially with increasing $k$ with $\lim_{a/k \to 0} \mathbf{E}[r(a,k)] \leq 2$. 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.

## 1   Introduction

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

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

**Instance** $<n, m, k>$**:** An universe $U = [1, n]$, sets $S_1, S_2, \ldots, S_m \subseteq U$ 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 $<n, m>$ instead of $<n, m, 1>$.

Both $\mathbf{SC}$ and $\mathbf{SC}_k$ are already well-known in the realm of design and analysis of combinatorial algorithms (*e.g.*, see [14]). 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**

**(a)** [4] *Assuming* $NP \not\subseteq DTIME(n^{\log \log n})$*, instances* $< n, m >$ *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)** [14] *An instance* $<n, m, k>$ *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 1 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 multicover problem. In Theorem 2 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 $\mathbf{e}^2 \approx 7.39$, and for

smaller values of $a/k$ it decreases towards 2 exponentially with increasing $k$ with $\lim_{a/k \to 0} \mathbf{E}[r(a,k)] \leq 2$. More precisely, $\mathbf{E}[r(a,k)]$ is at most

$1 + \ln a,$ if $k = 1$

$\left(1 + \mathbf{e}^{-(k-1)/5}\right) \ln(a/(k-1)),$ if $a/(k-1) \geq \mathbf{e}^2 \approx 7.39$ and $k > 1$

$\min\{\, 2 + 2 \cdot \mathbf{e}^{-(k-1)/5},\, 2 + \left(\mathbf{e}^{-2} + \mathbf{e}^{-9/8}\right) \cdot \frac{a}{k} \,\}$
$\approx \min\{\, 2 + 2 \cdot \mathbf{e}^{-(k-1)/5},\, 2 + 0.46 \cdot \frac{a}{k} \,\}$  if $a/(k-1) < \mathbf{e}^2$ and $k > 1$

*Some proofs are omitted due to lack of space.*

### 1.2  Summary of Analysis Techniques

– To prove Lemma 1, we generalize the approach in Johnson's paper [6]. 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 2 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 2 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, 2, 10, 14] were not always sufficient and hence we had to devise more appropriate bounds for certain parameter ranges.

## 2  Motivations

In this section is to 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 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 $B_1, B_2, \ldots, B_m$, can be retrieved with a *unit-cost query*;
– the columns of $B$ are in *general position*, *i.e.*, each subset of $k \leq n$ columns of $B$ is *linearly independent*;
– the *zero structure* of the matrix $C = AB = (c_{ij})$ is known, *i.e.*, a binary matrix $C^0 = \left(c_{ij}^0\right) \in \{0,1\}^{n \times m}$ is given, and it is known that $c_{ij} = 0$ for each $i,j$ for which $c_{ij}^0 = 0$.

The objective, "roughly speaking", 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).

Notice that there are 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, 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 $J = \{j_1, \ldots, j_\ell\}$. Then, the information obtained about $A$ may be summarized as $A_i \in \mathcal{H}_{J,i}^\perp$, where "$\perp$" indicates *orthogonal complement*, and

$$\mathcal{H}_{J,i} = \text{span}\ \{B_j, j \in J_i\}\ ,$$

$$J_i = \{j \mid j \in J \text{ and } c_{ij}^0 = 0\}. \tag{1}$$

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

$$\text{each set } J_i,\ i = 1, \ldots, n,\ \text{ has cardinality } \geq n - k, \tag{2}$$

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 most desirable special case 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 "$(+, +, -, 0, 0, -, \ldots)$") 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:

find $J$ of minimal cardinality such that $|J_i| \geq n - 1$, $i = 1, \ldots, n$. **(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.

More generally, suppose that the queries that we performed satisfy (2), 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 $\{v_1, \ldots, v_k\}$ 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 $\lambda_1, \ldots, \lambda_k$. 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 $v_r = \operatorname{col}(v_{1r}, \ldots, v_{nr})$, for each $r = 1, \ldots, k$. Since no coordinate of $A_i$ is zero, we know that $A_i$ belongs to the set $\mathcal{C} = \mathbb{R}^{k-1} \setminus \left( L_1 \bigcup \ldots \bigcup L_n \right)$ where, for each $1 \leq s \leq n$, $L_s$ is the hyperplane in $\mathbb{R}^{k-1}$ consisting of all those vectors $(\lambda_2, \ldots, \lambda_k)$ 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 [3, 11], and also [12] 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, $\Phi(n, d) = \sum_{i=0}^{d} \binom{n}{i} \leq 2\dfrac{n^d}{d!} \leq \left( \dfrac{\mathbf{e}n}{d} \right)^d$. 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 $i = 1, \ldots, n$.* Recalling the definition (1) of $J_i$, we see that $J_i = J \bigcap T_i$, 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

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

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

> *find $J \subseteq \{1, \ldots, m\}$ of minimal cardinality such that $|J \bigcap T_i| \geq n - k$,*
> $1 \leq i \leq n.$  **(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 $\{1, \ldots, n\}$ (so $m = n$); in this case, the set $J$ must equal $\{1, \ldots, m\}$ 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 $\{1, \ldots, n\}$ (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 problem was motivated by the setup for reverse-engineering of protein and gene networks described in [8, 9] and reviewed in [13]. We assume that the time evolution of a vector of state variables $x(t) = (x_1(t), \ldots, x_n(t))$ is described by a system of differential equations:

$$\dot{x}_1 = f_1(x_1, \ldots, x_n, p_1, \ldots, p_m)$$
$$\dot{x}_2 = f_2(x_1, \ldots, x_n, p_1, \ldots, p_m)$$
$$\vdots$$
$$\dot{x}_n = f_n(x_1, \ldots, x_n, p_1, \ldots, p_m)$$

(in vector form, "$\dot{x} = f(x, p)$"), where $p = (p_1, \ldots, p_m)$ is a vector of parameters, representing for instance the concentrations of certain enzymes which are maintained at a constant value during a particular experiment. There is a reference value $\bar{p}$ of $p$, which represents "wild type" (that is, normal) conditions, and a corresponding steady state $\bar{x}$. That is, $f(\bar{x}, \bar{p}) = 0$. 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 [8, 9], is that, while we do not know the form of $f$, we 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. This knowledge will be summarized by the binary matrix $C^0 = \left(c_{ij}^0\right) \in \{0, 1\}^{n \times m}$, 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 us to make small perturbations in which we change one of the parameters, say the $k$th one, while leaving the remaining ones constant. (A generalization would allow for the simultaneous perturbation of more than one parameter.) For the perturbed vector $p \approx \bar{p}$, we measure the resulting steady state vector $x = \xi(p)$. (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. In practice, each

such perturbation experiment involves letting the system relax to steady state, and the use of some biological reporting mechanism, such as microarrays, in order to measure the expression profile of the variables $x_i$.) For each of the possible $m$ experiments, in which a given $p_j$ is perturbed, we may estimate the $n$ "sensitivities"

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

(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 [8, 9] 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 $A = BC$. 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)**)[4]
**Instance $< m, n, k >$:** $U = [1, m]$ and sets $T_1, T_2, \ldots, T_n \subseteq U$ 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 $<m, n, k>$ of $\mathbf{CP}_k$ and an instance $<n, m, n - k>$ 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 $\mathbf{SC}_{n-k}$.

# 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 $c > 0$.

---

[4] $\mathbf{CP}_{n-1}$ is known as the hitting set problem [5, p. 222].

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

Johnson [6] provides an example in which the greedy heuristic for some instance of **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 1.** *For any fixed $c > 0$, the greedy heuristic (as described in Fact $1$(b)) has an approximation ratio of at least $\left(\frac{1}{2} - o(1)\right)\left(\frac{n-c}{8n-2}\right)\log_2 n = \Omega(\log n)$ for some instance $<n, m, n - c>$ of $\text{SC}_{n-c}$.*

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

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

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.* $\mathbf{1} = s$;
- $S^i = \{A \in \mathcal{S} : i \in A\}$ denotes the sets in $\mathcal{S}$ that contains a specific element $i$.

Consider the standard integer programming (IP) formulation of an instance $<n, m, k>$ of $\text{SC}_k$ [14]:

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

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 $\text{SC}_k$ can then be designed:

> **1.** Select an appropriate positive constant $\beta > 1$ in the following manner:
> $$\beta = \begin{cases} \ln a & \text{if } k = 1 \\ \ln(a/(k-1)) & \text{if } a/(k-1) \geq \mathbf{e}^2 \text{ and } k > 1 \\ 2 & \text{otherwise} \end{cases}$$
> **2.** Find a solution $x$ to the LP relaxation via any polynomial-time algorithm for solving linear programs (e.g. [7]).
> **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 $\mathbf{Pr}[A \in \mathcal{C}^1] = \beta x_A$.
> **5. (greedy selection)** Form a family of sets $\mathcal{C}^2$ as:
> while $s^i(c^0 + c^1 + c^2) < k$ for some $i \in U$, insert to $C^2$ any $A \in S^i - C^0 - C^1 - C^2$.
> **6.** Return $\mathcal{C} = \mathcal{C}^0 \cup \mathcal{C}^1 \cup \mathcal{C}^2$ as our solution.

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

**Theorem 2.**[5]

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

Let OPT denote the minimum number of sets used by an optimal solution. Obviously, OPT$\geq \mathbf{1}x$ and OPT$\geq \frac{nk}{a}$. A proof of Theorem 2 follows by showing the following upper bounds on $\mathbf{E}[r(a,k)]$ and taking the best of these bounds for each value of $a/k$:

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

**3.2.1   Proof of $\mathbf{E}[r(a,k)] \leq 1 + \ln a$ if $k = 1$,**
**$\mathbf{E}[r(a,k)] \leq \left(1 + \mathbf{e}^{-(k-1)/5}\right) \ln(a/(k-1))$ if $a/(k-1) \geq \mathbf{e}^2$ and $k > 1$,**
**and**
**$\mathbf{E}[r(a,k)] \leq 2 + 2 \cdot \mathbf{e}^{-(k-1)/5}$ if $a/(k-1) < \mathbf{e}^2$ and $k > 1$**

For our analysis, we first define two following two vector notations:

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

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 *bonus* $= \beta \mathbf{1}x^0 - \mathbf{1}c^0$. It is easy to see that $\mathbf{E}[\mathbf{1}(c^0 + c^1)] = \beta \mathbf{1}x - bonus$.

---

[5] The case of $k = 1$ was known before and included for the sake of completeness only.

The contribution of set $A$ to *bonus* is $\beta x_A^0 - c_A^0$. This contribution to *bonus* can be distributed equally to the elements if $A$. Since $|A| \le a$, an element $i \in [1, n]$ receives a total of *at least* $b^i/a$ of *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, $\mathbf{E}[r(a, k)] = \mathbf{E}[\mathbf{1}(c^0 + c^1 + c^2)] \le \beta \mathbf{1} x + \sum_{i=1}^n \mathbf{E}[y^i - \frac{b^i}{a}]$ 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 $b^i = s^i(\beta x^0 - c^0) \ge (\beta - 1)(k - g^i - q^i)$;
- $q^i \ge 0$, $g^i \ge 0$ and $g^i + q^i \le k$;
- $y^i$ 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 2.** $\mathbf{E}[Y_g^\mu] < \mathbf{e}^{-\mu} \sum_{i=0}^{g-1} \frac{g-i}{i!} \mu^i$.

From now on we will assume the worst-case distribution of $Y_g^\mu$, so we will assume that the above inequality in Lemma 2 is actually an equality (as it becomes so in the limit), *i.e.*, we assume $\mathbf{E}[Y_g^\mu] = \mathbf{e}^{-\mu} \sum_{i=0}^{g-1} \frac{g-i}{i!} \mu^i$. For a fixed $\beta$, we will need to estimate the growth of $\mathbf{E}[Y_g^{g\beta}]$ as a function of $g$. Let $\rho_g(\beta) = \mathbf{e}^{g\beta} \mathbf{E}[Y_g^{g\beta}]$.

**Lemma 3.** $\rho_g(1) = \sum_{i=0}^{g-1} \frac{g-i}{i!} g^i = \frac{g^g}{(g-1)!}$

**Lemma 4.** For $\beta > 1$, $\frac{\rho_{g+1}(\beta)}{\beta \rho_g(\beta)}$ is a decreasing function of $\beta$.

**Lemma 5.** If $g > 1$ and $\beta > 1$ then $\frac{\mathbf{E}[Y_g^{g\beta}]}{\mathbf{E}[Y_{g-1}^{(g-1)\beta}]} \le \mathbf{e}^{-\beta} \left(\frac{g}{g-1}\right)^g$

**Lemma 6.** $\frac{\mathbf{E}[Y_g^{g\beta+q}]}{\mathbf{E}[Y_g^{g\beta}]} < \mathbf{e}^{-q(1-1/\beta)}$

### 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 $\mathbf{E}[r(a, k)]$. 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 $y = 0$. We omit the case of $k = 1$ and assume that $k > 1$ for the rest of this section. 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 $p = \beta g = \beta$, $\mathbf{E}[y] = \mathbf{E}[Y_1^\beta] = \mathbf{e}^{-\beta}$.

We postulate that

$$
\begin{aligned}
\mathbf{E}[y - \frac{b}{a}] \leq 0 &\equiv \mathbf{e}^{-\beta} \leq \frac{(\beta - 1)(k - 1)}{a} \\
&\equiv \frac{\mathbf{e}^{-\beta}}{\beta - 1} \leq \frac{k - 1}{a} \\
&\equiv \mathbf{e}^{\beta}(\beta - 1) \geq \frac{a}{k - 1} \\
&\equiv \beta + \ln(\beta - 1) \geq \ln \frac{a}{k - 1} \quad (3)
\end{aligned}
$$

It is easy to see that, for the base case, $\mathbf{E}[\mathbf{1}(c^0 + c^1 + c^2)] \leq \beta \mathbf{1} x \leq \ln(a/(k - 1))$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 $(\beta - 1)(g + q - 1)/a$. Also, $\mathbf{E}[y] = \mathbf{E}[Y_g^p] = \mathbf{E}[Y_g^{\beta g + (\beta - 1)q}]$.

**Lemma 7.** $\frac{\mathbf{E}[Y_g^{\beta g + (\beta - 1)q}]}{\mathbf{E}[Y_1^\beta]} \leq \mathbf{e}^{-(g + q - 1)/5}$.

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 $\mathbf{e}^{-t/5}$. As a function of $t = g + q - 1$ we have

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

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 $\frac{(\beta - 1)(k - 1)}{a}\mathbf{e}^{-(k - 1)/5}$. As a result, our expected performance ratio for $k > 1$ is

given by

$$\mathbf{E}[r(a,k)] \leq \beta\mathbf{1}x + \sum_{i=1}^{n}\mathbf{E}[y^i - \tfrac{b^i}{a}]$$

$$\leq \beta\mathrm{OPT} + \tfrac{\beta nk}{a}\mathbf{e}^{-(k-1)/5}$$

$$\leq \beta(1 + \mathbf{e}^{-(k-1)/5})\mathrm{OPT}$$

$$\leq \begin{cases} \left(1 + \mathbf{e}^{-(k-1)/5}\right)\ln(a/(k-1))\,\mathrm{OPT} & \text{if } a/(k-1) \geq \mathbf{e}^2 \\ 2\cdot\left(1 + \mathbf{e}^{-(k-1)/5}\right)\mathrm{OPT} & \text{if } a/(k-1) < \mathbf{e}^2 \end{cases}$$

## References

1. N. Alon and J. Spencer, *The Probabilistic Method*, Wiley Interscience, New York, 1992.
2. 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.
3. 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.
4. U. Feige. *A threshold for approximating set cover*, JACM, Vol. 45, 1998, pp. 634-652.
5. M. R. Garey and D. S. Johnson. *Computers and Intractability - A Guide to the Theory of NP-Completeness*, W. H. Freeman & Co., 1979.
6. D. S. Johnson. *Approximation Algorithms for Combinatorial Problems*, Journal of Computer and Systems Sciences, Vol. 9, 1974, pp. 256-278.
7. N. Karmarkar. *A new polynomial-time algorithm for linear programming*, Combinatorica, 4: 373–395, 1984.
8. 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.
9. B. N. Kholodenko and E.D. Sontag, *Determination of functional network structure from local parameter dependence data*, arXiv physics/0205003, May 2002.
10. R. Motwani and P. Raghavan, *Randomized Algorithms*, Cambridge University Press, New York, NY, 1995.
11. L. Schläfli, *Theorie der vielfachen Kontinuitat (1852)*, in *Gesammelte Mathematische Abhandlungen*, volume 1, pp. 177–392, Birkhäuser, Basel, 1950.
12. 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.
13. 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.
14. V. Vazirani. *Approximation Algorithms*, Springer-Verlag, July 2001.