

The Minrank of Random Graphs

Alexander Golovnev^{*†}Oded Regev^{*†}Omri Weinstein^{*‡}

Abstract

The *minrank* of a graph G is the minimum rank of a matrix M that can be obtained from the adjacency matrix of G by switching ones to zeros (i.e., deleting edges) and setting all diagonal entries to one. This quantity is closely related to the fundamental information-theoretic problems of (linear) *index coding* (Bar-Yossef et al., FOCS'06), network coding and distributed storage, and to Valiant's approach for proving superlinear circuit lower bounds (Valiant, Boolean Function Complexity '92).

We prove tight bounds on the minrank of random Erdős-Rényi graphs $G(n, p)$ for all regimes of $p \in [0, 1]$. In particular, for any constant p , we show that $\text{minrk}(G) = \Theta(n/\log n)$ with high probability, where G is chosen from $G(n, p)$. This bound gives a near quadratic improvement over the previous best lower bound of $\Omega(\sqrt{n})$ (Haviv and Langberg, ISIT'12), and partially settles an open problem raised by Lubetzky and Stav (FOCS '07). Our lower bound matches the well-known upper bound obtained by the “clique covering” solution, and settles the linear index coding problem for random knowledge graphs.

Finally, our result suggests a new avenue of attack, via derandomization, on Valiant's approach for proving superlinear lower bounds for logarithmic-depth semilinear circuits.

1 Introduction

In information theory, the *index coding* problem [BK98, BYBJK06] is the following: A sender wishes to *broadcast* over a noiseless channel an n -symbol string $x \in \mathbb{F}^n$ to a group of n receivers R_1, \dots, R_n , each equipped with some *side information* (a subset $K_i \subseteq \{x_1, \dots, x_n\} \setminus \{x_i\}$). The index coding problem asks what is the minimum length m of a broadcast message that allows each receiver R_i to retrieve the i th symbol x_i , given his side-information K_i and the broadcasted message. The side information of the receivers is modeled by a directed graph \mathcal{K}_n with no self-loops, in which R_i observes $K_i := \{x_j : (i, j) \in E(\mathcal{K}_n)\}$. \mathcal{K}_n is sometimes called the *knowledge graph*. A canonical example is where \mathcal{K}_n is the complete graph on the vertex set $[n]$, i.e., each receiver observes all but his own symbol. In this simple case, broadcasting the sum $\sum_{i=1}^n x_i \pmod{\mathbb{F}}$ allows each receiver to retrieve his own symbol, hence $m = 1$.

This problem is motivated by applications to distributed storage [AK15], on-demand video streaming (ISCOD, [BK06]) and wireless networks (see, e.g., [YZ99]), where a typical scenario is that clients miss information during transmissions of the network, and the network is interested in minimizing the retransmission length by exploiting the side information clients already

^{*}Courant Institute of Mathematical Sciences, New York University.

[†]Supported by the Simons Collaboration on Algorithms and Geometry and by the National Science Foundation (NSF) under Grant No. CCF-1320188. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the NSF.

[‡]Supported by a Simons Foundation junior society fellowship.

possess. In theoretical computer science, index coding is related to some important communication models and problems in which players have overlapping information, such as the *one-way* communication complexity of the index function [KNR95] and the more general problem of *network coding* [ACLY00, ERL15]. Index coding can also be viewed as an interesting special case of nondeterministic computation in the (notoriously difficult to understand) multiparty *Number-On-Forehead* model, which in turn is a promising approach for proving data structure and circuit lower bounds [Pat10, PRS97, JS11]. The minimum length of an index code for a given graph has well-known relations to other important graph parameters – for an undirected graph G , it is bounded from below by the size of the maximum independent set ($\alpha(G)$), and for any graph (directed or not), it is bounded from above by the clique-cover number ($\chi(\bar{G})$) since for every clique in G , it suffices to broadcast a single symbol (recall the example above). The aforementioned connections also led to algorithmic connections (via convex relaxations) between the computational complexity of graph coloring and that of computing the minimum index code length of a graph [CH14].

In the context of circuit lower bounds, Riis [Rii07] observed that a certain index coding problem is equivalent to the so-called *shift conjecture* of Valiant [Val92] (see Subsection 1.1 below). If true, this conjecture would resolve a major open problem of proving superlinear size lower bound for logarithmic-depth circuits.

When the encoding function of the index code is *linear* (in x), the corresponding scheme is called a *linear index code*. Note that the example featured in the beginning of this introduction is actually a linear index coding scheme. In their seminal paper, Bar-Yossef et al. [BYBJK06] showed that the minimum length m of a *linear* index code is characterized precisely by a graph operator of the knowledge graph \mathcal{K}_n , called the *minrank* ($\text{minrk}_{\mathbb{F}}(\mathcal{K}_n)$), first introduced by Haemers [Hae79] in the context of Shannon capacity of graphs. Informally, $\text{minrk}_{\mathbb{F}}(\mathcal{K}_n)$ is the minimum rank (over \mathbb{F}) of an $n \times n$ matrix M that “represents” \mathcal{K}_n . By “represents” we mean a matrix M that can be obtained from the adjacency matrix of \mathcal{K}_n by *deleting edges* (i.e., turning non-zero entries to zeros) and setting all *diagonal* entries of M to non-zero elements (see Definition 1 for the formal definition). Note that without the “diagonal constraint”, the above minimum would trivially be 0, and indeed this constraint is what makes the problem interesting and hard to analyze. While linear index codes are in fact optimal for a large class of knowledge graphs (including directed acyclic graphs, perfect graphs, odd “holes” and odd “anti-holes” – complements of odd holes [BYBJK06]), there are examples where non-linear codes outperform their linear counterparts [LS07]. In the same paper, Lubetzky and Stav [LS07] posed the question about *typical* knowledge graphs, namely,

What is the minimum length of an index code for a random knowledge graph $\mathcal{K}_n = \mathcal{G}_{n,p}$?

Here, $\mathcal{G}_{n,p}$ denotes a random Erdős-Rényi directed graph, i.e., a graph on n vertices in which each arc is taken independently with probability p . In this paper, we partially answer this open problem by determining the optimal length of *linear* index codes for typical knowledge graphs. We prove a tight lower bound on the minrank of $\mathcal{G}_{n,p}$ for all values of $p \in [0, 1]$. In particular,

Theorem 1 (Main theorem, informal). *For any constant $0 < p < 1$ and any field \mathbb{F} of cardinality $|\mathbb{F}| < n^{O(1)}$, it holds with high probability that*

$$\text{minrk}_{\mathbb{F}}(\mathcal{G}_{n,p}) = \Theta\left(\frac{n}{\log n}\right).$$

The formal quantitative statement of our result can be found in Corollary 2 below. We note that our general result (see Theorem 2) extends beyond the constant regime to *subconstant* values

of p , and this feature of our lower bound is crucial for potential applications of our result to circuit lower bounds (we elaborate on this in the next subsection). Theorem 1 gives a near quadratic improvement over the previous best lower bound of $\Omega(\sqrt{n})$ [LS07, HL12], and settles the linear index coding problem for random knowledge graphs, as an $O_p(n/\log n)$ linear index coding scheme is achievable via the clique-covering solution (see Section 1.2).

In the following subsection, we propose a concrete (yet admittedly still quite challenging) approach for proving superlinear circuit lower bounds based on a potential “derandomization” of Theorem 1.

1.1 Connections to circuit lower bounds for semilinear circuits

In his seminal line of work, Valiant [Val77, Val83, Val92] proposed a path for proving superlinear lower bounds on the size of circuits with logarithmic depth. Informally speaking, Valiant’s “depth reduction” method [Val77, Vio09] allows, for any constant ε , to reduce any circuit of size $O(n)$ and depth $O(\log n)$ (with n inputs and n outputs), to a new circuit with the same inputs and outputs, where now each output gate is an (arbitrary) Boolean function of (i) at most n^ε inputs which are “hard-wired” to this output gate, and (ii) an additional fixed set of $m = O_\varepsilon(n/\log \log n)$ “common bits” $b_1(x), \dots, b_m(x)$ which in general may be arbitrary Boolean functions of the input $x = x_1, \dots, x_n$. Therefore, if one could exhibit a function that cannot be computed in this model using $O(n/\log \log n)$ common bits, this would imply a superlinear circuit lower bound for logarithmic depth circuits.

Valiant [Val92] proposed a concrete class of candidate hard functions for this new model. This candidate class of functions is the set of all cyclic shifts of an n -bit string x (i.e., the circuit must output $y = x \circ i$ which is the string obtained by shifting all the bits of x by i positions). Valiant conjectured that no “pre-wired” circuit as above can realize *all* n cyclic shifts using $mo(n/\log \log n)$ common bits (in fact, Valiant postulated that $m = \Omega(n)$ common bits are required, and this still seems plausible). This conjecture is sometimes referred to as *Valiant’s shift conjecture*.¹

As noted earlier in the introduction, Riis [Rii07] observed that a certain index coding problem is equivalent to this conjecture. Let $G = (V, A)$ be a directed graph, and $i \in \{0, \dots, n-1\}$. We denote by G^i the graph with vertex set V and arcs set $A^i = \{(u, v + i \pmod n) : (u, v) \in A\}$. Riis [Rii07] showed that the following conjecture is equivalent to Valiant’s shift conjecture:

Conjecture 1. *There exists $\varepsilon > 0$ such that for all sufficiently large n and every graph G on n vertices with out-degree of each vertex at most n^ε , there exists a shift i such that the minimum length of an index coding scheme for G^i (over \mathbb{F}_2) is $\omega(n/\log \log n)$.*

Naturally, one could consider a weaker version of Conjecture 1, in which the goal is to prove that there is some $i \in [n]$ for which the minimum length of a *linear* index-code for G^i must be large, namely, that $\text{minrk}_2(G^i) \geq \omega(n/\log \log n)$ for some $i \in [n]$. Let us refer to this as Conjecture 1b.

Let us consider a function $f(x, p)$ whose input is partitioned in two parts $x \in \{0, 1\}^k$ and $p \in \{0, 1\}^t$. We say that the function f is *semilinear* if for every fixed value of $p = p_0$, the function $f(x, p_0)$ is a linear function of the elements of x . The class of semilinear functions includes bilinear

¹It is noteworthy to mention that for the purpose of this conjecture, one could replace cyclic shifts with any (efficiently computable) subset of $\exp(O(n))$ permutations from S_n (indeed, since the permutation itself is part of the input in the above model, its description size must be linear in n , or else proving superlinear lower bounds would be trivial).

functions (such as matrix and integer multiplication), permutation p of the x part of the input and cyclic shifts. A circuit G is called *semilinear* if for every fixed value of $p = p_0$, one can assign linear functions to the gates of G , so that G computes $f(x, p_0)$.

It is easy to see that a semilinear function with one-bit output can always be computed by a linear-size log-depth circuit. However, if we consider semilinear functions with $O(n)$ output bits, then the circuit complexity of a random function is $\Omega(n^2/\log n)$ with high probability. It is an open problem to prove a superlinear lower bound against log-depth semilinear circuits [PRS97]. As Conjecture 1 is equivalent to Valiant’s shift conjecture, Conjecture 1b is equivalent to an analogous conjecture for semilinear circuits. Therefore, proving Conjecture 1b would in turn imply superlinear circuit lower bounds for logarithmic-depth *semilinear* circuits.

Theorem 1 suggests an approach for proving Conjecture 1b. Indeed, to prove Conjecture 1b, one needs to show that

$$\forall G \exists i \in [n] \text{ minrk}_2(G \circ i) \geq \Omega(n/\log \log n) ,$$

where G is an n -vertex directed graph with out-degree at most n^ϵ . Theorem 1 (and the more precise concentration bound we prove in Theorem 2) asserts that *almost all* n -vertex graphs of the “right” degree (i.e., $p = n^{\epsilon-1}$ for the expected degree of each vertex to be n^ϵ as in the above setting of the conjectures) have *linear* ($\Omega(n)$) minrank. Therefore, a conceivable approach is to try to “derandomize” Theorem 1, e.g., by showing that composing an arbitrary such knowledge graph G (possibly with low minrank) with a “pseudo-random” permutation (e.g., an $(n/\log n)$ -wise independent permutation, which can be described using only $O(n)$ bits [AL12]) leads to a high minrank.

It is worth noting that for this “derandomization” approach to work, it is crucial that Theorem 1 achieves the tight lower bound (up to constant factors) even for the sub-constant regime where $p = n^{\epsilon-1}$, in which case the bound shows that the minrank of $\mathcal{G}_{n,p}$ is $\Omega(n) \gg n/\log \log n$ (which is consistent with Valiant’s conjecture).

Finally, we mention one last circuit class for which the above “derandomization” approach might be easier. A circuit $G = (V, A)$ is called *Valiant series-parallel (VSP)*, if there is a labeling of its vertices $l: V \rightarrow \mathbb{R}$, such that for every arc $(u, v) \in A$, $l(u) < l(v)$, but there is no pair of arcs $(u, v), (u', v') \in A$, such that $l(u) < l(u') < l(v) < l(v')$. Most of the known circuit constructions (i.e., circuit upper bounds) are VSP circuits. Thus, it is also a big open question in circuit complexity to prove a superlinear lower bound on size of (semilinear) VSP circuit. Valiant [Val77] and Calabro [Cal08] show that for any *constant* $d \geq 1$, a series-parallel circuit of linear size can be reduced to a circuit where each output is a function of $O_d(1)$ input bits and $m = n/d$ common bits. Thus, derandomization of Theorem 1 for the regime of $p = O(n^{-1})$ could potentially lead to lower bounds for VSP circuits. Note that in the case of $p = O(n^{-1})$, the entropy of a random graph is only $O(n \log n)$ bits, hence, information-theoretically it seems easier to derandomize than the case of $p = n^{\epsilon-1}$.

1.2 Proof overview of Theorem 1

In [LS07], Lubetzky and Stav showed that for any field \mathbb{F} and a directed graph G ,

$$\text{minrk}_{\mathbb{F}}(G) \cdot \text{minrk}_{\mathbb{F}}(\bar{G}) \geq n .$$

This inequality gives a lower bound of $\Omega(\sqrt{n})$ on the expected value of the minrank of $\mathcal{G}_{n,1/2}$ (Indeed, the random variables $\mathcal{G}_{n,1/2}$ and $\bar{\mathcal{G}}_{n,1/2}$ have identical distributions). Since $\text{minrk}_{\mathbb{F}}(\mathcal{G}_{n,p})$

is monotonically non-increasing in p , the same bound holds for any $p \leq 1/2$. Note that this approach does not give bounds for $p > 1/2$, nor does it give strong concentration results. Haviv and Langberg [HL12] improved this result by proving a lower bound of $\Omega(\sqrt{n})$ for all constant p (that holds with high probability).

We now outline the main ideas of our proof. For simplicity we assume that $\mathbb{F} = \mathbb{F}_2$ and $p = 1/2$. To prove that $\text{minrk}_2(\mathcal{G}_{n,p}) \geq k$, it suffices to show that with high probability $\mathcal{G}_{n,p}$ has no representing matrix (in the sense of Definition 1) whose rank is less than k . A natural approach for doing so is trying to argue that any matrix with non-zero diagonal entries whose rank is $< k$, must be very *dense* (note again that without the diagonal assumption, this claim is completely false as we could choose the all-zeroes matrix and all bets are off). This would clearly be useful, as the probability that $\mathcal{G}_{n,p}$ is represented by any fixed dense representing matrix is exponentially small in the number of non-zero entries (arcs) the matrix contains. Indeed, a simple yet important observation in our proof, is that the sparsity s (the number of non-zero entries) of any rank- k matrix with 1s on its main diagonal must be at least $\approx n^2/k$.² This fact alone can already recover the previous $k \geq \Omega(\sqrt{n})$ lower bound of [HL12] for all constant p (albeit with a much weaker concentration bound). Indeed, since there are $\approx 2^{2nk}$ $n \times n$ matrices of rank k (as a rank- k matrix can be written as a product of $n \times k$ by $k \times n$ matrices, which requires $2nk$ bits to specify), the above fact implies, by a union bound, that the probability that $\mathcal{G}_{n,p}$ contains a subgraph of rank $\geq k$ is bounded from above by (roughly) $2^{2nk} \cdot (1/2)^{n^2/k}$, which is $\ll 1$ for $k = O(\sqrt{n})$.

In order to go beyond \sqrt{n} , we give a “compression” argument that implies a much tighter upper bound on the number of rank- k $n \times n$ matrices *with sparse bases*. To this end we show that every matrix is uniquely determined by specifying a row-basis, column-basis and the indices of these rows and columns in the matrix (see Lemma 2). Thus, we can encode a matrix using only $\approx s_{basis} \cdot \log n$ bits, where s_{basis} is the total number of ones in the sparsest column and row k -bases of the matrix. This in turn implies that there are only $\approx n^{s_{basis}}$ such matrices, which is potentially far less than the naive 2^{2nk} bound. Indeed, if we could argue that any rank- k matrix of sparsity s contains column and row bases of sparsity $O(k \cdot (s/n))$ (which is what one would expect as the average column sparsity is s/n), then we would be done: The union bound would now yield (roughly)

$$2^{ks \log(n)/n} \cdot (1/2)^s,$$

so one could set the minrank parameter k to be as large as $\Theta(n/\log n)$ and the above expression would still be $\ll 1$. Hence to complete the proof, we need to show that every matrix has column/row bases whose (per-column/row) sparsity matches the average sparsity of the columns (rows). Perhaps surprisingly, this turns out to be the most challenging part of the proof.

Indeed, it is not hard to see that there are matrices for which such a sparse basis does not exist (see the simple example discussed before Lemma 4). Nevertheless, we show that every matrix must contain a large³ *submatrix* with such sparse basis relative to its rank (here we mean a “proper” square submatrix of possibly lower dimension and rank, but we show that ruling out these submatrices suffices). At a very high level, in order to find a sparse-basis submatrix, we sort the

²To see why, notice that any maximal linearly independent set of columns must “cover” all coordinates, i.e., there must not be any coordinate that is zero in all vectors, as otherwise we could take the column vector corresponding to that coordinate and it would be linearly independent of our set (due to the nonzero diagonal) in contradiction to maximality. Assuming all columns have roughly the same number of 1s, we obtain that each column has at least n/k 1s, leading to the claimed bound. See Lemma 3 for the full proof.

³In fact, for this lemma to be useful, the right measure of “largeness” turns out to be the *ratio* between the rank and the dimension of the sub-matrix, see Lemma 4 for the precise statement.

columns in an increasing order of their sparsities, then greedily construct a column basis vector by vector. While the argument works for any matrix, we show that if the matrix has non-zero entries on the diagonal, then we necessarily reach a point where the average sparsity of the currently constructed basis does not exceed the average sparsity of submatrix it spans, hence yielding a “nontrivial” submatrix. See Lemma 4 for the full proof.

2 Preliminaries

For an integer n , we denote the set $\{1, \dots, n\}$ by $[n]$. For an integer n and $0 \leq p \leq 1$, we denote by $\mathcal{G}_{n,p}$ the probability space over the directed graphs on n vertices where each arc is taken independently with probability p .

For a directed graph G , by $\chi(G)$ we denote the chromatic number of undirected graph that has the same set of vertices as G , and an edge in place of every arc of G . By \bar{G} we mean a directed graph on the same set of vertices as G that contains an arc if and only if G does not contain it.⁴

Let \mathbb{F} be a finite field. For a vector $v \in \mathbb{F}^n$, we denote by v^j the j th entry of v , and by $v^{\leq j} \in \mathbb{F}^j$ the vector v truncated to its first j coordinates. For a matrix $M \in \mathbb{F}^{n \times n}$ and indices $i, j \in [n]$, let $M_{i,j}$ be the entry in the i th row and j th column of M , $\text{Col}_i(M)$ be the i th column of M , $\text{Row}_i(M)$ be the i th row of M , and $\text{rk}(M)$ be the rank of M over \mathbb{F} .

By a *principal submatrix* we mean a submatrix whose set of row indices is the same as the set of column indices. By the *leading principal submatrix* of size k we mean a principal submatrix that contains the first k columns and rows.

For a matrix $M \in \mathbb{F}^{n \times n}$, the sparsity $s(M)$ is the number of non-zero entries in M . We say that a matrix $M \in \mathbb{F}^{n \times n}$ of rank k *contains* an *s -sparse column (row) basis*, if M *contains* a column (row) basis (i.e., a set of k linearly independent columns (rows)) with a total of at most s non-zero entries.

Definition 1 (Minrank [BYBJK06, LS07]).⁵ *Let $G = (V, A)$ be a graph on $n = |V|$ vertices with the set of directed arcs A . A matrix $M \in \mathbb{F}^{n \times n}$ represents G if $M_{i,i} \neq 0$ for every $i \in [n]$, and $M_{i,j} = 0$ whenever $(i, j) \notin A$ and $i \neq j$. The minrank of G over \mathbb{F} is*

$$\text{minrk}_{\mathbb{F}}(G) = \min_{M \text{ represents } G} \text{rk}(M).$$

We say that two graphs *differ at only one vertex* if they differ only in arcs leaving one vertex. Following [HHMS10, HL12], to amplify the probability in Theorem 2, we shall use the following form of Azuma’s inequality for the vertex exposure martingale.

Lemma 1 (Corollary 7.2.2 and Theorem 7.2.3 in [AS16]). *Let $f(\cdot)$ be a function that maps directed graphs to \mathbb{R} . If f satisfies the inequality $|f(H) - f(H')| \leq 1$ whenever the graphs H and H' differ at only one vertex, then*

$$\Pr[|f(\mathcal{G}_{n,p}) - \mathbb{E}[f(\mathcal{G}_{n,p})]| > \lambda\sqrt{n-1}] < 2e^{-\lambda^2/2}.$$

⁴Throughout the paper we assume that graphs under consideration do not contain self-loops. In particular, neither G nor \bar{G} has self-loops.

⁵In this paper we consider the directed version of minrank. Since the minrank of a directed graph does not exceed the minrank of its undirected counterpart, a lower bound for a directed random graph implies the same lower bound for an undirected random graph. The bound is tight for both directed and undirected random graphs (see Theorem 3).

3 The Minrank of a Random Graph

The following elementary linear-algebraic lemma shows that a matrix $M \in \mathbb{F}^{n \times n}$ of rank k is fully specified by k linearly independent rows, k linearly independent columns, and their $2k$ indices. In what follows, we denote by $\mathcal{M}_{n,k}$ the set of matrices from $\mathbb{F}^{n \times n}$ of rank k .

Lemma 2 (Row and column bases encode the entire matrix). *Let $M \in \mathcal{M}_{n,k}$, and let $R = (\text{Row}_{i_1}(M), \dots, \text{Row}_{i_k}(M))$, $C = (\text{Col}_{j_1}(M), \dots, \text{Col}_{j_k}(M))$ be, respectively, a row basis and a column basis of M . Then the mapping $\phi: \mathcal{M}_{n,k} \rightarrow (\mathbb{F}^{1 \times n})^k \times (\mathbb{F}^{n \times 1})^k \times [n]^{2k}$ defined as*

$$\phi(M) = (R, C, i_1, \dots, i_k, j_1, \dots, j_k),$$

is a one-to-one mapping.

Proof. We first claim that the intersection of R and C has full rank, i.e., that the submatrix $M' \in \mathbb{F}^{k \times k}$ obtained by taking rows i_1, \dots, i_k and columns j_1, \dots, j_k has rank k . This is a standard fact, see, e.g., [HJ13, p20, Section 0.7.6]. We include a proof for completeness. Assume for convenience that $(i_1, \dots, i_k) = (1, \dots, k)$ and $(j_1, \dots, j_k) = (1, \dots, k)$. Next, assume towards contradiction that $\text{rk}(M') = \text{rk}(\{\text{Col}_1(M'), \dots, \text{Col}_k(M')\}) = k' < k$. Since C is a column basis of M , every column $\text{Col}_i(M)$ is a linear combination of vectors from C , and in particular, every $\text{Col}_i(M')$ is a linear combination of $\{\text{Col}_1(M'), \dots, \text{Col}_k(M')\}$. Therefore, the $k \times n$ submatrix $M'' := (\text{Col}_1^{\leq k}(M), \dots, \text{Col}_n^{\leq k}(M))$ has rank k' . On the other hand, the k rows of M'' : $\text{Row}_1(M), \dots, \text{Row}_k(M)$ were chosen to be linearly independent by construction. Thus, $\text{rk}(M'') = k > k'$, which leads to a contradiction.

In order to show that ϕ is one-to-one, we show that R and C (together with their indices) uniquely determine the remaining entries of M . We again assume for convenience that $(i_1, \dots, i_k) = (1, \dots, k)$ and $(j_1, \dots, j_k) = (1, \dots, k)$. Consider any column vector $\text{Col}_i(M)$, $i \in [n] \setminus [k]$. By definition, $\text{Col}_i(M) = \sum_{t=1}^k \alpha_{i,t} \cdot \text{Col}_t(M)$ for some coefficient vector $\alpha_i := (\alpha_{i,1}, \dots, \alpha_{i,k}) \in \mathbb{F}^{k \times 1}$. Thus, in order to completely specify all the entries of $\text{Col}_i(M)$, it suffices to determine the coefficient vector α_i . But M' has full rank, hence the equation

$$M' \alpha_i^T = \text{Col}_i^{\leq k}(M)$$

has a *unique* solution. Therefore, the coefficient vector α_i is fully determined by M' and $\text{Col}_i^{\leq k}(M)$. Thus, the matrix M can be uniquely recovered from R, C and the indices $\{i_1, \dots, i_k\}, \{j_1, \dots, j_k\}$. \square

The following corollary gives us an upper bound on the number of low-rank matrices that contain sparse column and row bases. In what follows, we denote by $\mathcal{M}_{n,k,s}$ the set of matrices over $\mathbb{F}^{n \times n}$ of rank k that contain an s -sparse row basis and an s -sparse column basis.

Corollary 1 (Efficient encoding of sparse-base matrices).

$$|\mathcal{M}_{n,k,s}| \leq (n \cdot |\mathbb{F}|)^{6s}.$$

Proof. Throughout the proof, we assume without loss of generality that $s \geq k$, as otherwise $|\mathcal{M}_{n,k,s}| = 0$ hence the inequality trivially holds. The function ϕ from Lemma 2 maps matrices from $\mathcal{M}_{n,k,s}$ to $(R, C, i_1, \dots, i_k, j_1, \dots, j_k)$, where R and C are s -sparse bases. Therefore, the total number of matrices in $\mathcal{M}_{n,k,s}$ is bounded from above by

$$\left(\binom{kn}{s} \cdot |\mathbb{F}|^s \right)^2 \cdot n^{2k} \leq ((n^2)^s \cdot |\mathbb{F}|^s)^2 \cdot n^{2k} \leq (n \cdot |\mathbb{F}|)^{6s},$$

where the last inequality follows from $k \leq s$. □

Now we show that a matrix of low rank with nonzero entries on the main diagonal must contain many nonzero entries. To get some intuition on this, notice that a rank 1 matrix with nonzero entries on the diagonal must be nonzero everywhere. Also notice that the assumption on the diagonal is crucial – low rank matrices in general can be very sparse.

Lemma 3 (Sparsity vs. Rank for matrices with non-zero diagonal). *For any matrix $M \in \mathbb{F}^{n \times n}$ with non-zero entries on the main diagonal (i.e., $M_{i,i} \neq 0$ for all $i \in [n]$), it holds that*

$$s(M) \geq \frac{n^2}{4\text{rk}(M)}.$$

Proof. Let s denote $s(M)$. The average number of nonzero entries in a column of M is s/n . Therefore, Markov's inequality implies that there are at least $n/2$ columns in M each of which has sparsity at most $2s/n$. Assume without loss of generality that these are the first $n/2$ columns of M . Now pick a maximal set of linearly independent columns among these columns. We claim that the cardinality of this set is at least $n^2/(4s)$. Indeed, in any set of less than $n^2/(4s)$ columns, the number of coordinates that are nonzero in any of the columns is less than

$$\frac{n^2}{4s} \cdot \frac{2s}{n} = \frac{n}{2}$$

and therefore there exists a coordinate $i \in \{1, \dots, n/2\}$ that is zero in all those columns. As a result, the i th column, which by assumption has a nonzero i th coordinate, must be linearly independent of all those columns, in contradiction to the maximality of the set. We therefore get that

$$\text{rk}(M) \geq n^2/(4s),$$

as desired. □

The last lemma we need is also the least trivial. In order to use Corollary 1, we would like to show that any $n \times n$ matrix of rank k has sparse row and column bases, where by sparse we mean that their sparsity is roughly k/n times that of the entire matrix. If the number of nonzero entries in each row and column was roughly the same, then this would be trivial, as we can take any maximal set of linearly independent columns or rows. However, in general, this might be impossible to achieve. E.g., consider the $n \times n$ matrix whose first k columns are chosen uniformly and the remaining $n - k$ columns are all zero. Then any column basis would have to contain all first k columns (since they are linearly independent with high probability) and hence its sparsity is equal to that of the entire matrix. Instead, what the lemma shows is that one can always choose a *principal submatrix* with the desired property, i.e., that it contains sparse row and column bases, while at the same time having relative rank that is at most that of the original matrix.

Lemma 4 (Every matrix contains a principal submatrix of low relative-rank and sparse bases). *Let $M \in \mathcal{M}_{n,k}$ be a matrix. There exists a principal submatrix $M' \in \mathcal{M}_{n',k'}$ of M , such that $k'/n' \leq k/n$, and M' contains a column basis and a row basis of sparsity at most*

$$s(M') \cdot \frac{2k'}{n'}.$$

Note that if M contains a zero entry on the main diagonal, the lemma becomes trivial. Indeed, we can take M' to be a 1×1 principal submatrix formed by this zero entry. Thus, the lemma is only interesting for matrices M without zero elements on the main diagonal (i.e., when every principal submatrix has rank greater than 0).

Proof. We prove the statement of the lemma by induction on n . The base case $n = 1$ holds trivially.

Now let $n > 1$, and assume that the statement of the lemma is proven for every $m \times m$ matrix for $1 \leq m < n$. Let $s(i)$ be the number of nonzero entries in the i th column plus the number of non-zero entries in the i th row (note that a nonzero entry on the diagonal is counted twice). Let also $s_{\max} = \max_i s(i)$. By applying the same permutation to the columns and rows of M we can assume that $s(1) \leq s(2) \leq \dots \leq s(n)$ holds.

If for some $1 \leq n' < n$, the leading principal submatrix M' of dimensions $n' \times n'$ has rank at most $k' \leq n'k/n$, then we use the induction hypothesis for M' . This gives us a principal submatrix M'' of dimensions $n'' \times n''$ and rank k'' , such that M'' contains a column basis and a row basis of sparsity at most $s(M'') \cdot \frac{2k''}{n''}$. Also, by induction hypothesis $k''/n'' \leq k'/n' \leq k/n$, which proves the lemma statement in this case.

Now we assume that for all $n' < n$, the rank of the leading principal submatrix of dimension $n' \times n'$ is greater than $n'k/n$. We prove that the lemma statement holds for $M' = M$ for a column basis, and an analogous proof gives the same result for a row basis.

For every $0 \leq i \leq s_{\max}$, let $a_i = |\{j : s(j) = i\}|$. Note that

$$\sum_{i=0}^{s_{\max}} a_i = n. \quad (1)$$

Let us select a column basis of cardinality k by greedily adding linearly independent vectors to the basis in non-decreasing order of $s(i)$. Let k_i be the number of selected vectors j with $s(j) = i$. Then

$$\sum_{i=0}^{s_{\max}} k_i = k. \quad (2)$$

Next, for any $0 \leq t < s_{\max}$, consider the leading principal submatrix given by indices i with $s(i) \leq t$. The rank of this matrix is at most $k' = \sum_{i=0}^t k_i$, and its dimensions are $n' \times n'$, where $n' = \sum_{i=0}^t a_i < n$. Thus by our assumption $k'/n' \geq k/n$, or equivalently,

$$\sum_{i=0}^t k_i \geq \frac{k}{n} \cdot \sum_{i=0}^t a_i. \quad (3)$$

From (1) and (2),

$$\sum_{i=0}^{s_{\max}} k_i = \frac{k}{n} \cdot \sum_{i=0}^{s_{\max}} a_i. \quad (4)$$

Now, (3) and (4) imply that for all $0 \leq t \leq s_{\max}$:

$$\sum_{i=t}^{s_{\max}} k_i \leq \frac{k}{n} \cdot \sum_{i=t}^{s_{\max}} a_i. \quad (5)$$

To finish the proof, notice that the sparsity of the constructed basis of M is at most

$$\sum_{i=1}^{s_{\max}} i \cdot k_i = \sum_{t=1}^{s_{\max}} \sum_{i=t}^{s_{\max}} k_i \stackrel{(5)}{\leq} \frac{k}{n} \cdot \sum_{t=1}^{s_{\max}} \sum_{i=t}^{s_{\max}} a_i = \frac{k}{n} \cdot \sum_{i=1}^{s_{\max}} i \cdot a_i = s(M) \cdot \frac{2k}{n}.$$

□

Now we are ready to prove our main result – a lower bound on the minrank of a random graph.

Theorem 2.

$$\Pr \left[\text{minrk}_{\mathbb{F}}(\mathcal{G}_{n,p}) \geq \Omega \left(\frac{n \log(1/p)}{\log(n|\mathbb{F}|/p)} \right) \right] \geq 1 - e^{-\Omega \left(\frac{n \log^2(1/p)}{\log^2(n|\mathbb{F}|/p)} \right)}.$$

Proof. Let us bound from above probability that a random graph $\mathcal{G}_{n,p}$ has minrank at most

$$k := \frac{n \log(1/p)}{C \log(n|\mathbb{F}|/p)},$$

for some constant C to be chosen below.

Recall that by Lemma 4, every matrix of rank at most k contains a principal submatrix $M' \in \mathcal{M}_{n',k'}$ of sparsity $s' = s(M')$ with column and row bases of sparsity at most

$$s' \cdot \frac{2k}{n},$$

where $k'/n' \leq k/n$. By Corollary 1, there are at most $(n' \cdot |\mathbb{F}|)^{6(2s'k/n)}$ such matrices M' , and (for any s') there are $\binom{n}{n'}$ ways to choose a principal submatrix of size n' in a matrix of size $n \times n$. Furthermore, recall that Lemma 3 asserts that for every n', k' ,

$$s' \geq \frac{n'^2}{4k'}. \quad (6)$$

Finally, since M' contains at least $s' - n'$ off-diagonal non-zero entries, $\mathcal{G}_{n,p}$ contains it with probability at most $p^{s'-n'}$. We therefore have

$$\begin{aligned} \Pr [\text{minrk}_{\mathbb{F}}(\mathcal{G}_{n,p}) \leq k] &\leq \sum_{k',n',s'} \Pr \left[\mathcal{G}_{n,p} \text{ contains } M' \in \mathcal{M}_{n',k'}, s(M') = s', s(\text{bases of } M') \leq s' \cdot \frac{2k}{n} \right] \\ &\leq \sum_{k',n',s'} \binom{n}{n'} \cdot p^{s'-n'} \cdot (n' \cdot |\mathbb{F}|)^{12s'k/n} \\ &\leq \sum_{k',n',s'} 2^{n' \log n - s' \log(1/p) + n' \log(1/p) + (12s'k/n) \log(n'|\mathbb{F}|)}, \end{aligned} \quad (7)$$

where all the summations are taken over n', k' , s.t. $k'/n' \leq k/n$ and $s' \geq \frac{n'^2}{4k'}$, and the first inequality is again by Lemma 4. We now argue that for sufficiently large constant C , all positive terms in the exponent of (7) are dominated by the magnitude of the negative term ($s' \log(1/p)$). Indeed:

$$\begin{aligned} n' \log n + n' \log(1/p) + (12s'k/n) \log(n'|\mathbb{F}|) &= n' \log(n/p) + (12s'k/n) \log(n'|\mathbb{F}|) \\ &\leq (4s'k'/n') \log(n/p) + (12s'k/n) \log(n|\mathbb{F}|) \leq (16s'k/n) \log(n|\mathbb{F}|/p) = (16s'/C) \log(1/p), \end{aligned}$$

where the first inequality follows from (6), and the second one follows from $k'/n' \leq k/n$.

Thus, for $C > 16$,

$$\Pr \left[\text{minrk}_{\mathbb{F}}(\mathcal{G}_{n,p}) \leq \frac{n \log(1/p)}{C \log(n|\mathbb{F}|/p)} \right] \leq n^4 \cdot 2^{-\Omega(s' \log(1/p))} \leq 2^{-\Omega(\log(n))}.$$

In particular, $\mathbb{E}[\text{minrk}_{\mathbb{F}}(\mathcal{G}_{n,p})] \geq \frac{n \log(1/p)}{2C \log(n|\mathbb{F}|/p)}$. Furthermore, note that changing a single row (or column) of a matrix can change its minrank by at most 1, hence the minrank of two graphs that differ in one vertex differs by at most 1. We may thus apply Lemma 1 with $\lambda = \Theta\left(\frac{\sqrt{n} \log(1/p)}{\log(n|\mathbb{F}|/p)}\right)$ to obtain

$$\Pr \left[\text{minrk}_{\mathbb{F}}(\mathcal{G}_{n,p}) \geq \Omega\left(\frac{n \log(1/p)}{\log(n|\mathbb{F}|/p)}\right) \right] \geq 1 - e^{-\Omega\left(\frac{n \log^2(1/p)}{\log^2(n|\mathbb{F}|/p)}\right)}.$$

as desired. \square

Corollary 2. *For a constant $0 < p < 1$ and a field \mathbb{F} of size $|\mathbb{F}| < n^{O(1)}$,*

$$\Pr[\text{minrk}_{\mathbb{F}}(\mathcal{G}_{n,p}) \geq \Omega(n/\log n)] \geq 1 - e^{-\Omega(n/\log^2 n)}.$$

3.1 Tightness of Theorem 2

In this section, we show that Theorem 2 provides a tight bound for all values of p bounded away from 1 (i.e., $p \leq 1 - \Omega(1)$). (See also the end of the section for the regime of p close to 1.)

Theorem 3. *For any p bounded away from 1,*

$$\Pr \left[\text{minrk}_{\mathbb{F}}(\mathcal{G}_{n,p}) = O\left(\frac{n \log(1/p)}{\log n + \log(1/p)}\right) \right] \geq 1 - e^{-\Omega(n)}.$$

Proof. We can assume that $p > n^{-1/8}$ as otherwise the statement is trivial.

As we saw in the introduction, in the case of a clique (a graph with an arc between every pair of distinct vertices) it is enough to broadcast only one bit. This simple observation leads to the ‘‘clique-covering’’ upper bound: If a directed graph G can be covered by m cliques, then $\text{minrk}_{\mathbb{F}}(G) \leq m$ [Hae78, BYBJK06, HL12]. Note that the minimal number of cliques needed to cover G is exactly $\chi(\bar{G})$. Thus, we have the following upper bound: For any field \mathbb{F} and any directed graph G ,

$$\text{minrk}_{\mathbb{F}}(G) \leq \chi(\bar{G}). \tag{8}$$

Since the complement of $\mathcal{G}_{n,p}$ is $\mathcal{G}_{n,1-p}$, it follows from (8) that an upper bound on $\chi(\mathcal{G}_{n,1-p})$ implies an upper bound on $\text{minrk}_{\mathbb{F}}(\mathcal{G}_{n,p})$.

Let $\mathcal{G}_{n,p}^-$ denote a random Erdős-Rényi *undirected* graph on n vertices, where each edge is drawn independently with probability p . For constant p , the classical result of Bollobás [Bol88] asserts that the chromatic number of an undirected random graph satisfies

$$\Pr \left[\chi(\mathcal{G}_{n,1-p}^-) \leq \frac{n \log(1/p)}{2 \log n} (1 + o(1)) \right] > 1 - e^{-\Omega(n)}. \tag{9}$$

In fact, Pudlák, Rödl, and Sgall [PRS97] showed that (9) holds for any $p > n^{-1/4}$.

Since we define the chromatic number of a directed graph to be the chromatic number of its undirected counterpart, $\chi(\mathcal{G}_{n,1-p}) = \chi(\mathcal{G}_{n,1-p}^-)$. The bound (9) depends on p only logarithmically ($\log(1/p)$), thus, asymptotically the same bounds hold for the chromatic number of a random directed graph. \square

The lower bound of Theorem 2 is also almost tight for the other extreme regime of $p = 1 - \varepsilon$, where $\varepsilon = o(1)$. Łuczak [Luc91] proved that for $p = 1 - \Omega(1/n)$,

$$\Pr \left[\chi(\mathcal{G}_{n,1-p}^-) \leq \frac{n(1-p)}{2 \log n(1-p)} (1 + o(1)) \right] > 1 - (n(1-p))^{-\Omega(1)}. \quad (10)$$

When $p = 1 - \varepsilon$, the upper bound (10) matches the lower bound of Theorem 2 for $\varepsilon \geq n^{-1+\Omega(1)}$. For $\varepsilon = O(n^{-1})$, (10) gives an asymptotically tight upper bound of $O(1)$. Thus, we only have a gap between the lower bound of Theorem 2 and known upper bounds when $p = 1 - \varepsilon$ and $\omega(1) \leq n\varepsilon \leq n^{o(1)}$.

Acknowledgements

We would like to thank Ishay Haviv for his valuable comments on an earlier version of this work.

References

- [ACLY00] Rudolf Ahlswede, Ning Cai, Shuo-Yen Robert Li, and Raymond W. Yeung. Network information flow. *IEEE Transactions on information theory*, 46(4):1204–1216, 2000.
- [AK15] Fatemeh Arbabjolfaei and Young-Han Kim. Three stories on a two-sided coin: Index coding, locally recoverable distributed storage, and guessing games on graphs. *CoRR*, abs/1511.01050, 2015.
- [AL12] Noga Alon and Shachar Lovett. Almost k -wise vs. k -wise independent permutations, and uniformity for general group actions. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques - 15th International Workshop, APPROX 2012, and 16th International Workshop, RANDOM 2012, Cambridge, MA, USA, August 15-17, 2012. Proceedings*, pages 350–361, 2012.
- [AS16] N. Alon and J.H. Spencer. *The Probabilistic Method*. Wiley Series in Discrete Mathematics and Optimization. Wiley, 2016.
- [BK98] Yitzhak Birk and Tomer Kol. Informed-source coding-on-demand (ISCOD) over broadcast channels. In *Proceedings IEEE INFOCOM '98, The Conference on Computer Communications, Seventeenth Annual Joint Conference of the IEEE Computer and Communications Societies, Gateway to the 21st Century, San Francisco, CA, USA, March 29 - April 2, 1998*, pages 1257–1264, 1998.
- [BK06] Yitzhak Birk and Tomer Kol. Coding on demand by an informed source (ISCOD) for efficient broadcast of different supplemental data to caching clients. *IEEE Trans. Information Theory*, 52(6):2825–2830, 2006.

- [Bol88] Béla Bollobás. The chromatic number of random graphs. *Combinatorica*, 8(1):49–55, 1988.
- [BYBJK06] Ziv Bar-Yossef, Yitzhak Birk, TS Jayram, and Tomer Kol. Index coding with side information. In *Foundations of Computer Science, 2006. FOCS'06. 47th Annual IEEE Symposium on*, pages 197–206. IEEE, 2006.
- [Cal08] Chris Calabro. A lower bound on the size of series-parallel graphs dense in long paths. In *Electronic Colloquium on Computational Complexity (ECCC)*, volume 15, 2008.
- [CH14] Eden Chlamtac and Ishay Haviv. Linear index coding via semidefinite programming. *Combinatorics, Probability & Computing*, 23(2):223–247, 2014.
- [ERL15] Michelle Effros, Salim Y. El Rouayheb, and Michael Langberg. An equivalence between network coding and index coding. *IEEE Trans. Information Theory*, 61(5):2478–2487, 2015.
- [Hae78] Willem Haemers. An upper bound for the Shannon capacity of a graph. In *Colloq. Math. Soc. János Bolyai*, volume 25, pages 267–272, 1978.
- [Hae79] Willem Haemers. On some problems of Lovász concerning the Shannon capacity of a graph. *IEEE Transactions on Information Theory*, 25(2):231–232, 1979.
- [HHMS10] H. Tracy Hall, Leslie Hogben, Ryan Martin, and Bryan Shader. Expected values of parameters associated with the minimum rank of a graph. *Linear Algebra and its Applications*, 433(1):101–117, 2010.
- [HJ13] Roger A. Horn and Charles R. Johnson. *Matrix analysis*. Cambridge University Press, Cambridge, second edition, 2013.
- [HL12] Ishay Haviv and Michael Langberg. On linear index coding for random graphs. In *Information Theory Proceedings (ISIT), 2012 IEEE International Symposium on*, pages 2231–2235. IEEE, 2012.
- [JS11] Stasys Jukna and Georg Schnitger. Min-rank conjecture for log-depth circuits. *J. Comput. Syst. Sci.*, 77(6):1023–1038, 2011.
- [KNR95] Ilan Kremer, Noam Nisan, and Dana Ron. On randomized one-round communication complexity. In *Proceedings of the Twenty-seventh Annual ACM Symposium on Theory of Computing, STOC '95*, pages 596–605, New York, NY, USA, 1995. ACM.
- [LS07] Eyal Lubetzky and Uri Stav. Non-linear index coding outperforming the linear optimum. In *Foundations of Computer Science, 2007. FOCS'07. 48th Annual IEEE Symposium on*, pages 161–168. IEEE, 2007.
- [Luc91] Tomasz Łuczak. The chromatic number of random graphs. *Combinatorica*, 11(1):45–54, 1991.
- [Pat10] Mihai Patrascu. Towards polynomial lower bounds for dynamic problems. In *Proceedings of the 42nd ACM Symposium on Theory of Computing, STOC 2010, Cambridge, Massachusetts, USA, 5-8 June 2010*, pages 603–610, 2010.

- [PRS97] Pavel Pudlák, Vojtech Rödl, and Jirí Sgall. Boolean circuits, tensor ranks, and communication complexity. *SIAM Journal on Computing*, 26(3):605–633, 1997.
- [Rii07] Søren Riis. Information flows, graphs and their guessing numbers. *Electr. J. Comb.*, 14(1), 2007.
- [Val77] Leslie G. Valiant. Graph-theoretic arguments in low-level complexity. In *Mathematical Foundations of Computer Science 1977, 6th Symposium, Tatranska Lomnica, Czechoslovakia, September 5-9, 1977, Proceedings*, pages 162–176, 1977.
- [Val83] Leslie G. Valiant. Exponential lower bounds for restricted monotone circuits. In *Proceedings of the fifteenth annual ACM symposium on Theory of computing*, pages 110–117. ACM, 1983.
- [Val92] Leslie G. Valiant. Why is Boolean complexity theory difficult. *Boolean Function Complexity*, 169:84–94, 1992.
- [Vio09] Emanuele Viola. On the power of small-depth computation. *Foundations and Trends in Theoretical Computer Science*, 5(1):1–72, 2009.
- [YZ99] Raymond W. Yeung and Zhen Zhang. Distributed source coding for satellite communications. *IEEE Trans. Information Theory*, 45(4):1111–1120, 1999.