

Rank of the vertex-edge incidence matrix of r -out hypergraphs

Colin Cooper

Department of Informatics
King's College
London WC2B 4BG
England

Alan Frieze*

Department of Mathematical Sciences
Carnegie Mellon University
Pittsburgh PA15213
U.S.A.

June 4, 2022

Abstract

We consider the rank of a class of sparse Boolean matrices of size $n \times n$. In particular, we show that the probability that such a matrix has full rank, and is thus invertible, is a positive constant with value about 0.2574 for large n .

The matrices arise as the vertex-edge incidence matrix of 1-out 3-uniform hypergraphs. The result that the null space is bounded in expectation, can be contrasted with results for the usual models of sparse Boolean matrices, based on the vertex-edge incidence matrix of random k -uniform hypergraphs. For this latter model, the expected co-rank is linear in the number of vertices n , [5], [8].

For fields of higher order, the co-rank is typically Poisson distributed.

General notes for the reviewers

1. Corrections and requests for comments raised by the reviewers (and other corrections) are shown in red text.
2. The labelling of Lemma 11 onwards has changed, as Remark 11 was redundant and has been removed.
3. Lemma 6 has been simplified as suggested.

*Research supported in part by NSF Grant DMS1661063

4. Lemma 14 (Lemma 15 as was). The notation has been reduced, corrected and made consistent. The 4 cases in the proof use the corrected notation throughout. Corrections to notation within the cases are not highlighted in red except for technical errors, or to improve explanation.

1 Introduction

For positive integers $r \geq 1$, $s \geq 2$, let $\mathbf{M}(s, r, n)$ be the space of $n \times rn$ matrices with entries generated in the following manner. For each $i = 1, \dots, n$ there are r columns $C_{i,j}$, $j = 1, \dots, r$. Each column $C_{i,j}$ has a unit entry in row i , and $s-1$ other unit entries, in rows chosen randomly with replacement from $[n]$, or without replacement from $[n] - \{i\}$, all other entries in the column being zero. In general we consider the arithmetic on entries in the matrix, (and thus the evaluation of linear dependencies), to be over $GF(2)$. If so, in the “with replacement case”, if two unit entries coincide the entry is set to zero. When $r = 1$, the matrix consists of an identity matrix plus $s-1$ random units in each column. If $s = 2$, and entries (and columns $C_{i,j}$, $j = 1, \dots, r$) are chosen without replacement, $\mathbf{M}(2, r, n)$ is the space of vertex-edge incidence matrices of the random graphs $G_{r\text{-out}}(n)$.

For fixed integer k , $\vec{G}_{k\text{-out}}(n)$ is a random digraph with vertex set $[n]$. The k arcs from any vertex v have terminal vertices chosen uniformly as any of the $\binom{n-1}{k}$ random k -subsets of $[n] \setminus \{v\}$. The multi-graph $G_{k\text{-out}}(n)$ is obtained from $\vec{G}_{k\text{-out}}(n)$ by ignoring the orientation of the edges. The $G_{k\text{-out}}(n)$ model of random graphs has been extensively studied, see e.g., Chapter 16 of [12] for an introduction. It is known to be k -connected for $k \geq 2$, Fenner and Frieze [10], to have a perfect matching for $k \geq 2$, Frieze [11], and to be Hamiltonian for $k \geq 3$, Bohman and Frieze [4].

If $s \geq 3$, then $M \in \mathbf{M}(s, r, n)$ is the vertex-edge incidence matrix of a random r -out, s -uniform hypergraph. Random Boolean matrices based on the vertex-edge incidence matrix of s -uniform hypergraphs where the columns (edges) are chosen i.i.d. from all columns with s ones were studied by Cooper, Frieze and Pegden, [8]. A very general paper by Coja-Oghlan, Ergür, Gao, Hetterich and Rolvien, [5], gives the limiting rank in this latter model for a wide range of assumptions on the distribution of non-zero entries in the rows and columns. The fundamental difference between the r -out model of random matrices, and those of [5], [8] is the presence of an $n \times n$ identity matrix as a sub-matrix (in the without replacement case).

We will use ρ to denote the (row) rank of our matrices and then the co-rank is $n - \rho$. If the field is $GF(2)$, $\mathbf{x} \in \{0, 1\}^n$ is a linear dependency (*dependency* for short) if $\mathbf{x}M = \mathbf{0}$. Let $|\mathbf{x}| = |\{j : x_j = 1\}|$. We say that a set of rows $D \subseteq [n]$ is a dependency if $D = \{j : x_j = 1\}$ for some dependency \mathbf{x} . An ℓ -dependency is one where $|\mathbf{x}| = \ell$ or $|D| = \ell$. Two sizes of

dependency occur frequently in our proofs. For brevity we will say a dependency \mathbf{x} is *small* if $|\mathbf{x}| \leq \omega$ where $\omega \rightarrow \infty$ slowly, and a dependency \mathbf{x} is *large* if $|\mathbf{x}| = n/2 + O(\sqrt{n \log n})$.

Of particular interest is the case $r = 1$ which gives $n \times n$ Boolean matrices. The space $\mathbf{M}(2, 1, n)$ corresponds to random functional digraphs. The co-rank of these matrices over $GF(2)$ is well understood, see e.g., [2], [12], so the first extant case is $r = 1, s = 3$. We will show that over $GF(2)$, for $r = 1, s = 3$, the linear dependencies among the rows of M are w.h.p. either *small* or *large*, and the distributions of these dependencies are somewhat entangled. Estimating the interaction between small and large dependencies in matrices from $\mathbf{M}(3, 1, n)$ is the main problem we solve.

For $r = 1, s = 3$, define a Poisson parameter ϕ for small dependencies. The value of ϕ differs between the “with replacement” ϕ_R , and “without replacement” models $\phi_{\bar{R}}$ as follows:

$$\phi_R = \sum_{\ell \geq 1} \frac{1}{\ell} (2e^{-2})^\ell \sum_{j=0}^{\ell-1} \frac{\ell^j}{j!}, \quad \phi_{\bar{R}} = \sum_{\ell \geq 2} \frac{1}{\ell} (2e^{-2})^\ell \sum_{j=0}^{\ell-2} \frac{\ell^j}{j!}. \quad (1)$$

The numeric values are $\phi_R \approx 0.5215$, and $\phi_{\bar{R}} \approx 0.1151$, where $a \approx b$ means approximately equal.

Let

$$P(\sigma, \lambda) = \left(\frac{1}{2}\right)^{\lambda(\lambda+\sigma)} \frac{1}{\prod_{j=1}^{\lambda} \left(1 - \left(\frac{1}{2}\right)^j\right)} \prod_{j=1}^{\infty} \left(1 - \left(\frac{1}{2}\right)^{\lambda+\sigma+j}\right). \quad (2)$$

The quantity $P(\sigma, \lambda)$ is the limiting value of $\mathbb{P}(\lambda \mid \sigma)$ of the conditional probability of $\lambda = d - \sigma$ given σ , where σ is the dimension of the space induced by *small* dependencies and d the dimension of the space induced by all dependencies.

Theorem 1. *Let the matrix M be chosen u.a.r. from $\mathbf{M}(3, 1, n)$. Let $d \geq 0$ be integer. Over $GF(2)$, the limiting probability that M has co-rank d is given by*

$$\lim_{n \rightarrow \infty} \mathbb{P}(\text{co-rank}(M) = d) = e^{-\phi} \sum_{\sigma=0}^d \frac{\phi^\sigma}{\sigma!} P(\sigma, d - \sigma). \quad (3)$$

In particular,

$$\mathbb{P}(\text{rank}(M) = n) \sim e^{-\phi} P(0, 0) = e^{-\phi} \prod_{j=1}^{\infty} \left(1 - \left(\frac{1}{2}\right)^j\right).$$

Theorem 1 differs from many previous results on sparse random Boolean matrices. The co-rank (dimension of the null space) is bounded in expectation, and the matrix is invertible

with probability $e^{-\phi}P(0, 0) \approx 0.2574$ in the without replacement model. The bounded co-rank given by Theorem 1 can be contrasted with results for the edge-vertex incidence matrix of random hypergraphs, ([5], [8]), where the expected co-rank is linear in the number of vertices n , and the probability of a full rank matrix is exponentially small.

The matrices $\mathbf{M}(3, 1, n)$ exhibit a gap in the size of the dependencies (small or large), which we next explain.

Theorem 2. *Let M be chosen u.a.r. from $\mathbf{M}(3, 1, n)$, then w.h.p. either (i) a dependency \mathbf{x} is small i.e. $|\mathbf{x}| \leq \omega$ where $\omega \rightarrow \infty$ slowly or (ii) \mathbf{x} is large i.e. $|\mathbf{x}| = n/2 + O(\sqrt{n \log n})$.*

A gap property in solutions to random XOR-SAT systems over $GF(2)$ was previously observed by Achlioptas and Molloy [1], and by Ibrahimi, Kanoria, Kranić and Montanari [13]. They found that the Hamming distance between XOR-SAT solutions was either $O(\log n)$ or at least αn ; where n is the number of variables. In our case, large dependencies have intersection about $n/4$ (see Section 4), giving a precise value of α .

A dependency \mathbf{x} is *fundamental* if there is no other dependency $\mathbf{y} \neq \mathbf{x}$ such that $\mathbf{y} \leq \mathbf{x}$, componentwise. We will prove in Section 2 that the number Z of fundamental small dependencies is asymptotically distributed as $Po(\phi)$ i.e. Poisson with mean ϕ . The quantity $P(\sigma, \lambda)$ in (3) is the limiting probability that small dependencies span a space of dimension σ , and large dependencies increase the co-rank by λ .

Let π be the probability distribution given by

$$\pi(k) = \begin{cases} \prod_{j=1}^{\infty} \left(1 - \left(\frac{1}{2}\right)^j\right) & k = 0. \\ \frac{\prod_{j=k+1}^{\infty} \left(1 - \left(\frac{1}{2}\right)^j\right)}{\prod_{j=1}^k \left(1 - \left(\frac{1}{2}\right)^j\right)} \left(\frac{1}{2}\right)^{k^2} & k \geq 1. \end{cases} \quad (4)$$

Note that $\pi(k) = P(0, k)$ as given in (2). The probability distribution defined by π was previously observed in a model of random matrices over $GF(2)$ in which the entries $m_{i,j}$ are i.i.d. Bernoulli random variables with $\mathbb{P}(m_{i,j} = 1) = p$. For a wide range of p the distribution of dimension k of the null space is given by $\pi(k)$. The result was proved by Kovalenko et al., [14] for $p = 1/2$, and extended to the range $\min(p(n), 1 - p(n)) \geq (\log n + c(n))/n$, (where $c(n) \rightarrow \infty$ slowly) by Cooper [6]. A similar distributional result holds for the model of random matrices over the finite field $GF(q)$, see Cooper [7]. Here the non-zero entries $\alpha \in GF(q) \setminus \{0\}$ are independently and uniformly distributed with $\mathbb{P}(m_{i,j} = \alpha) = p/(q-1)$. The distribution of co-rank $\pi_q(k)$ equivalent to $\pi(k) = \pi_2(k)$ in (4) is obtained by replacing the $(1/2)$ terms in (4) by $(1/q)$.

Finally we mention some related cases for r -out s -uniform hypergraphs. For $r = 1$ and $s = 2$, M has expected rank $\sim n - (\log n)/2$. This is because the expected number of components in a random mapping is $\sim (1/2) \log n$, (see e.g., [12]). Note: For s even, the rows of M add to zero modulo 2. The following theorem is immediate from the proof of Theorem 1.

Theorem 3. *If $r \geq 2$ and $s = 2, 3$, then M has rank $n^* = n - \mathbb{1}_{\{s=2\}}$, w.h.p.*

The proof of Theorem 3, and results for finite fields of character $q \geq 3$ can be found in [9].

Notation: Apart from $O(\cdot), o(\cdot), \Omega(\cdot)$ as a function of $n \rightarrow \infty$, we use the notation $A_n \sim B_n$ if $\lim_{n \rightarrow \infty} A_n/B_n = 1$. The symbol $a \approx b$ indicates approximate numerical equality due to decimal truncation. The notation $\omega(n)$ describes a function tending to infinity as $n \rightarrow \infty$. The expression *with high probability* (w.h.p.), means with probability $1 - o(1)$, where the $o(1)$ is a function of n , which tends to zero as $n \rightarrow \infty$.

Outline of the proof for $GF(2)$ with $r = 1, s = 3$

Because the proofs are rather technical, *we give a detailed proof in the “with replacement” model*. For brevity, we omit the proof that the results are also valid in the “without replacement” model in this paper; the proof can be found in [9].

We refer to the rows of M as $M_i, i \in [n]$ and to the columns as $C_j, j \in [n]$. By a set of rows S , we mean the set of rows $M_i, i \in S$. A set of rows with indices L is linearly dependent (zero-sum) if $\sum_{i \in L} M_i = \mathbf{0} = \mathbf{0}(\text{mod } 2)$. A linear dependence L is *small* if $|L| \leq \omega$, where $\omega = \omega(n)$ is a function tending slowly to infinity with n . A linear dependence L is *large* if $|L| = (n/2)(1 + O(\sqrt{\log n/n}))$. As part of our proof, we show that w.h.p. there are no other sizes of dependency. A set of zero-sum rows L is *fundamental*, if L contains no smaller zero-sum set and L is disjoint from all other zero-sum sets. **This will be the case for minimal small dependencies, whereas** zero-sum sets of size about $n/2$ are not disjoint. We count k -sequences of large dependencies with a property we call *simple*. Many of the problems with the proofs arise because large dependencies are not disjoint, and are conditioned by the simultaneous presence of small dependencies in M .

We next outline the main steps in the proof of Theorem 1.

1. In Section 2 we prove that the number Z of small fundamental dependencies has factorial moments $\mathbf{E}(Z)_k \sim \phi^k$, where ϕ is given by (1). Thus Z is asymptotically Poisson distributed and

$$\mathbb{P}(M \text{ has } i \text{ small fundamental linear dependencies}) \sim \frac{\phi^i}{i!} e^{-\phi}.$$

2. For $M \in \mathbf{M}(3, 1, n)$ w.h.p. any fundamental sets of zero-sum rows of M are either small (of size $\ell \leq \omega$) or large (of size $\ell = (n/2)(1 + O(\sqrt{\log n/n}))$). This is proved in Section 3.
3. In Section 5 we discuss *simple* sequences of large dependencies, and in Section 6 we estimate the moments of these sequences and determine their interaction with small dependencies.

4. In Section 7 we estimate the number of simple sequences, conditional on the the number of small fundamental dependencies. This leads to an approximate set of linear equations whose solution completes the proof of Theorem 1.

2 Small dependencies in $GF(2)$: with replacement

Notation For $1 \leq k \leq \omega$, where $\omega \rightarrow \infty$ arbitrarily slowly with n , let $X_k(M)$ or $Y_k(M)$ denote the number of index sets of k -dependencies in M . A k -dependency is *small* if $k \leq \omega$. To distinguish the cases, we use Y_k when $k \leq \omega$, and use X_k when $k > \omega$. We will show that for values of $k > \omega$ other than $k \sim n/2$, $X_k = 0$ w.h.p. We use Z to denote the number of small fundamental dependent sets among the rows of M .

We first consider dependencies with $s = o(n^{1/2})$ rows. For $S \subseteq [n]$, let $\mathcal{F}(S)$ denote the event that the rows corresponding to S are dependent. Let Y_s denote the number of s -set dependencies.

Lemma 4. *If $|S| = s = o(n^{1/2})$ then*

$$\mathbb{P}(\mathcal{F}(S)) \sim \left(\frac{2s}{n}\right)^s e^{-2s}. \quad (5)$$

If $\omega \rightarrow \infty$, $\omega \leq s = o(n^{1/2})$ then $Y_s = 0$ w.h.p.

Proof. Suppose that $s = o(n^{1/2})$ and $S = [s]$. Then,

$$\begin{aligned} \mathbb{P}(\mathcal{F}(S)) &= \left(2\left(\frac{s}{n}\right)\left(\frac{n-s}{n}\right)\right)^s \left(\left(\frac{s}{n}\right)^2 + \left(\frac{n-s}{n}\right)^2\right)^{n-s} \\ &\sim \left(\frac{2s}{n}\right)^s e^{-2s}, \quad \text{using } s = o(\sqrt{n}). \end{aligned} \quad (6)$$

Explanation: The probability that exactly one of the two random choices in a column of S lies in a row of S is $2\left(\frac{s}{n}\right)\left(\frac{n-s}{n}\right)$. The probability that both or neither of the two random choices in a column of $[n] \setminus S$ lies in a row of S is $\left(\frac{s}{n}\right)^2 + \left(\frac{n-s}{n}\right)^2$.

This verifies (5). It follows that

$$\mathbf{E} Y_s \sim \binom{n}{s} \left(\frac{2s}{n}\right)^s e^{-2s} \sim \frac{(2s)^s e^{-2s}}{s!},$$

As $\mathbf{E} Y_{s+1}/\mathbf{E} Y_s \sim 2/e$ we have that $\mathbf{E} Y_\omega = e^{-\Omega(\omega)}$ and so w.h.p. there are no dependencies with $\omega \leq s = o(n^{1/2})$. \square

Define σ_s , κ_s by

$$\sigma_s = \sum_{j=0}^{s-1} \frac{s^j}{j!}, \quad \text{and} \quad \kappa_s = \frac{(s-1)!}{s^s} \sigma_s. \quad (7)$$

For $S \subseteq [n]$, let $\mathcal{F}^*(S)$ denote the event that the rows corresponding to S form a fundamental dependency. The next three lemmas deal with small fundamental dependencies.

Lemma 5. $\mathbb{P}(\mathcal{F}^*(S) \mid \mathcal{F}(S)) = \kappa_s$.

Proof. With high probability the rows S of a small dependency have the following structure: suppose that $|S| = s$. There is an $s \times s$ sub-matrix $M_{S,S}$ with unit diagonal entries and one random entry per column, and a zero $(s \times n-s)$ sub-matrix. For $i \in S$, either $M_{i,i} = 1$, and there is a unique entry $M_{j,i} = 1$ which gives rise to an *edge* (i, j) , or the random entry falls in position i in which case $M_{i,i} = 0$ and we regard this as a *loop* (i, i) . Thus $M_{S,S}$ is the incidence matrix of a random functional digraph D_S , and S is fundamental iff the underlying graph of D_S is connected. For $s \geq 1$, $\mathbb{P}(D_S \text{ is connected}) = \kappa_s$ (see e.g., [2], Theorem 14.33 or [12], Theorem 15.5). \square

Lemma 6. *Small fundamental dependent sets of M are pairwise disjoint, w.h.p.*

Proof. Let S, T be distinct fundamental zero-sum row sets with non-trivial intersection $S \cap T$. As functional digraphs have out-degree one, it follows that some column of $K = S \cup T$ must have three non-zero entries in the rows of K . Provided $\omega = o(\log n)$, the probability of such an event for $|K| \leq \omega$ is at most

$$\sum_{k=2}^{\omega} \binom{n}{k} \sum_{i=1}^k \binom{k}{i} \left(\frac{k}{n}\right)^{2i} \left(\frac{2k}{n}\right)^{k-i} = \frac{O(k^4)}{n} (2e)^k = o(1).$$

\square

Given this lemma we can now prove a Poisson distribution for Z .

Lemma 7. *The number Z of small fundamental dependent sets among the rows of M is asymptotically Poisson distributed with parameter ϕ_R , and thus*

$$\mathbb{P}(Z = d) \sim \frac{\phi_R^d}{d!} e^{-\phi_R}. \quad (8)$$

Proof. Fix $S \subseteq [n]$ and let S_1, \dots, S_d be a partition of S with $|S_i| = s_i, i = 1, 2, \dots, d$. Let $P(s_1, \dots, s_d)$ be the probability that each $S_i, i = 1, 2, \dots, d$ is a fundamental set, given that S is a dependency. Thus,

$$P(s_1, \dots, s_d) = \frac{(s_1)^{s_1} \cdots (s_d)^{s_d}}{s^s} \prod_{i=1, \dots, d} \mathbb{P}(D_{S_i} \text{ connected}) = \frac{1}{s^s} \prod_{i=1}^d (s_i - 1)! \sigma_{s_i}.$$

Explanation: the factor $\frac{(s_1)^{s_1} \cdots (s_d)^{s_d}}{s^s}$ is the conditional probability that the random choices for columns with index in S_i are in rows with index in S_i .

Thus, using (5), we see that

$$\mathbf{E}(Z)_d \sim \sum_{s \geq 1} \frac{(2s)^s}{s!} e^{-2s} \sum_{s_1 + \cdots + s_d = s} \binom{s}{s_1, \dots, s_d} P(s_1, \dots, s_d) \quad (9)$$

$$\begin{aligned} &= \sum_{s \geq 1} \sum_{s_1 + \cdots + s_d = s} \prod_{i=1}^d (2e^{-2})^{s_i} \frac{1}{s_i} \sigma_{s_i} \\ &= \left(\sum_{s \geq 1} \frac{1}{s} (2e^{-2})^s \sigma_s \right)^d \\ &= \phi_R^d. \end{aligned} \quad (10)$$

Thus, by the method of moments, the number of small disjoint fundamental zero-sum sets Z tends to a Poisson distribution with parameter ϕ_R . \square

3 Large zero-sum sets: First moment calculations

Define an index set J_a as follows,

$$J_a = \{n/2 - \sqrt{an \log n} \leq \ell \leq n/2 + \sqrt{an \log n}\} \text{ and } \bar{J}_a = [n] \setminus J_a, a \geq 0. \quad (11)$$

Lemma 8. (Large linearly dependent sets.) *Let X_ℓ denote the number of ℓ -dependencies among the rows of M .*

$$(i) \sum_{\ell \in J_1} \mathbf{E} X_\ell \sim 1.$$

$$(ii) \text{ Let } F = [n] \setminus ([\omega] \cup J_1), \text{ where } \omega \rightarrow \infty \text{ arbitrarily slowly with } n. \text{ Then } \sum_{\ell \in F} \mathbf{E} X_\ell = o(1).$$

Proof. From (6), the expected number of dependencies of size ℓ is

$$\mathbf{E} X_\ell = \binom{n}{\ell} \left(2 \left(\frac{\ell}{n} \right) \left(\frac{n-\ell}{n} \right) \right)^\ell \left(\left(\frac{\ell}{n} \right)^2 + \left(\frac{n-\ell}{n} \right)^2 \right)^{n-\ell}.$$

We next approximate the expression for $\mathbf{E} X_\ell$. We note the following expansion.

$$(1+x) \log(1-x^2) + (1-x) \log(1+x^2) = -2 \left(x^3 + \frac{x^4}{2} + \frac{x^7}{3} + \sum_{k \geq 4} \mathbb{1}_{\{k \text{ even}\}} \frac{x^{2k}}{k} \left(1 + \frac{kx^3}{k+1} \right) \right). \quad (12)$$

We write $\mathbf{E} X_\ell = \binom{n}{\ell} \Phi_\ell^n$, $\ell = (n/2)(1 + \varepsilon)$, where

$$\begin{aligned}
\Phi_\ell &= \left(\frac{1 - \varepsilon^2}{2} \right)^{\frac{(1+\varepsilon)}{2}} \left(\left(\frac{1 + \varepsilon}{2} \right)^2 + \left(\frac{1 - \varepsilon}{2} \right)^2 \right)^{\frac{(1-\varepsilon)}{2}} \\
&= \frac{1}{2} (1 - \varepsilon^2)^{\frac{(1+\varepsilon)}{2}} (1 + \varepsilon^2)^{\frac{(1-\varepsilon)}{2}} \\
&= \frac{1}{2} \exp \left\{ \frac{1}{2} \left((1 + \varepsilon) \log(1 - \varepsilon^2) + (1 - \varepsilon) \log(1 + \varepsilon^2) \right) \right\} \\
&= \frac{1}{2} \exp \left\{ - \left(\varepsilon^3 + \frac{\varepsilon^4}{2} + \frac{\varepsilon^7}{3} + \sum_{k \geq 4} \mathbb{1}_{\{k \text{ even}\}} \varepsilon^{2k} \left(\frac{1}{k} + \frac{\varepsilon^3}{k+1} \right) \right) \right\} \\
&= \frac{1}{2} \exp \left\{ - \left(\varepsilon^3 + \frac{\varepsilon^4}{2} + \varepsilon_7 \right) \right\}, \tag{13}
\end{aligned}$$

where $|\varepsilon_7| \leq 2|\varepsilon|^7/3$ for sufficiently small ε .

Also for $\ell = (n/2)(1 + \varepsilon)$, $|\varepsilon| < 1$,

$$\binom{n}{\ell} = \left(1 + O \left(\frac{1}{n} \right) \right) \frac{2^n}{\sqrt{2\pi n(1 - \varepsilon^2)}} \exp \left(-n \left(\frac{\varepsilon^2}{2} + \frac{\varepsilon^4}{12} + \varepsilon_6 \right) \right), \tag{14}$$

where $|\varepsilon_6| \leq |\varepsilon|^6/10$.

Case 1: $\ell \in J_1$. From (14) with $|\varepsilon| = 2\sqrt{(\log n)/n}$ we have

$$\frac{1}{2^n} \sum_{\ell \notin J_1} \binom{n}{\ell} = O(1/n^{5/2}),$$

so that

$$\frac{1}{2^n} \sum_{\ell \in J_1} \binom{n}{\ell} = 1 - O(1/n^{5/2}).$$

Using (13), for $\ell \in J_1$, $\Phi_\ell^n = e^{\Theta(n\varepsilon^3)}/2^n$. Then, as $n\varepsilon^3 = O(\log^{3/2} n/\sqrt{n})$,

$$\sum_{\ell \in J_1} \mathbf{E} X_\ell = \sum_{\ell \in J_1} \binom{n}{\ell} \frac{1}{2^n} e^{\Theta(n\varepsilon^3)} = 1 + o(1).$$

For future reference, we note that for $|\varepsilon| < c < 1$,

$$\begin{aligned}
\mathbf{E} X_\ell &= \binom{n}{\ell} \frac{1}{2^n} \exp \left\{ -n \left(\varepsilon^3 + \frac{\varepsilon^4}{2} + \varepsilon_7 \right) \right\} \\
&= \frac{(1+o(1))}{\sqrt{2\pi n(1-\varepsilon^2)}} \exp \left\{ -n \left(\frac{\varepsilon^2}{2} + \varepsilon^3 + \frac{\varepsilon^4}{2} + \frac{\varepsilon^4}{12} + \varepsilon_6 + \varepsilon_7 \right) \right\} \\
&= \frac{(1+o(1))}{\sqrt{2\pi n(1-\varepsilon^2)}} \exp \left\{ -\frac{n\varepsilon^2}{2} \left((1+\varepsilon)^2 + \frac{\varepsilon^2}{6} + O(\varepsilon^4) \right) \right\}. \tag{15}
\end{aligned}$$

Case 2: $\ell \in F$. Write $F = [n] \setminus ([\omega] \cup J_1)$ as $F = F_1 \cup F_2 \cup F_3$ where $F_1 = \{\omega, \dots, 3n/10\}$, $F_2 = \{7n/10, \dots, n\}$ and $F_3 = F \setminus (F_1 \cup F_2)$. Thus, for $\ell \in F_3$, $\ell = (n/2)(1+\varepsilon)$ where $-2/5 \leq \varepsilon \leq -\sqrt{(2 \log n)/n}$ or $\sqrt{(2 \log n)/n} \leq \varepsilon \leq 2/5$.

Case $\ell \in F_1$. For sufficiently large n , Stirling's approximation implies that

$$\binom{n}{\ell} \leq \frac{n^n}{\ell^\ell (n-\ell)^{n-\ell}},$$

so for some constant C (in both with and without replacement models)

$$\mathbf{E} X_\ell \leq \frac{Cn^n}{\ell^\ell (n-\ell)^{n-\ell}} \left(2 \left(\frac{\ell}{n} \right) \left(\frac{n-\ell}{n} \right) \right)^\ell \left(\left(\frac{\ell}{n} \right)^2 + \left(\frac{n-\ell}{n} \right)^2 \right)^{n-\ell}. \tag{16}$$

Continuing with this expression, using $\ell = \lambda n$ for $\lambda < 1/2$,

$$\begin{aligned}
\mathbf{E} X_\ell &\leq C \left(\frac{2^\lambda}{\lambda^\lambda (1-\lambda)^{1-\lambda}} \lambda^\lambda (1-\lambda)^\lambda (\lambda^2 + (1-\lambda)^2)^{1-\lambda} \right)^n \\
&= C \left(2^\lambda (1-\lambda)^\lambda \left(1 - \lambda + \frac{\lambda^2}{1-\lambda} \right)^{1-\lambda} \right)^n \\
&\leq C \left(2^\lambda (1-\lambda)^\lambda e^{-\lambda(1-\lambda)+\lambda^2} \right)^n \\
&= C (2(1-\lambda) e^{-1+2\lambda})^{\lambda n} \\
&= C [g(\lambda)]^{\lambda n}.
\end{aligned}$$

The function $g(\lambda)$ is strictly concave and has a unique maximum at $\lambda = 1/2$ with $g(1/2) = 1$. For $\lambda \leq 3/10$, $g(\lambda) \leq g(3/10) = (7/5)e^{-2/5} < 1$ so that

$$\sum_{\ell \in F_1} \mathbf{E} X_\ell \leq C \sum_{\ell \in F_1} g(3/10)^\ell = o(1).$$

Case $\ell \in F_2$. Referring to (15), the function $h(\varepsilon) = (\varepsilon^2/2)((1+\varepsilon)^2 + \varepsilon^2/6 + \varepsilon_6 + \varepsilon_7)$ satisfies $h(\varepsilon) > 2/25$ for $\varepsilon \geq 2/5$, and so

$$\sum_{\ell \in F_2} \mathbf{E} X_\ell \leq \sum_{\ell \in F_2} e^{-\Omega(n)} = o(1).$$

Case $\ell \in F_3$. For $\sqrt{(2 \log n)/n} \leq |\varepsilon| \leq \sqrt{(25 \log n)/n}$, the function $h(\varepsilon) \geq (1-o(1))(\log n)/n$. Let F_{3a} be the values of ℓ in this range

$$\sum_{\ell \in F_{3a}} \mathbf{E} X_\ell = O(\sqrt{n \log n}/n^{1-o(1)}) = o(1/n^{1/3}).$$

Let $F_{3b} = F_3 \setminus F_{3a}$. Then $\varepsilon^2/2 \geq (25/2)(\log n)/n$, and $(1+\varepsilon)^2 + \varepsilon^2/6 + \varepsilon_6 + \varepsilon_7 > 9/25$. Referring to (15),

$$\sum_{\ell \in F_{3b}} \mathbf{E} X_\ell = O(n)/n^4 = o(1/n^3).$$

□

4 Higher moments of large zero-sum sets: Background

Let $A \oplus B$ denote the symmetric set difference of the sets A and B . Thus $A \oplus B = (A \cup B) \setminus (A \cap B) = (A \setminus B) \cup (B \setminus A)$. Suppose that, over $GF(2)$, the rows $\textcolor{red}{M}_i, i \in A$ indexed by A are zero-sum, thus $\mathbf{z}_A = \sum_{i \in A} M[i] = \mathbf{0}$. Let B be another set such that $\mathbf{z}_B = \mathbf{0}$. We can write $\mathbf{z}_A = \mathbf{z}_{A \setminus B} + \mathbf{z}_{A \cap B}$ and $\mathbf{z}_B = \mathbf{z}_{B \setminus A} + \mathbf{z}_{A \cap B}$. Adding these two sets of rows modulo 2 has the effect of canceling the intersection $A \cap B$. Thus (i) $\mathbf{z}_A + \mathbf{z}_B = \mathbf{0}$, whether $\mathbf{z}_{A \cap B}$ is itself zero-sum or not; and (ii) $\mathbf{z}_A + \mathbf{z}_B = \mathbf{z}_{A \oplus B}$.

Recall that a set of zero-sum rows is fundamental if it contains no smaller zero-sum set of rows. For small sets we were able to count fundamental dependencies directly. We have to adopt an alternative strategy for large zero-sum sets. We use an approach similar to the one given in [6]. We count *simple* sequences of large linearly dependent row sets $B = (B_1, \dots, B_k)$, $k \geq 1$ constant, and where $|B_i| \in J_1$ so that $|B_i| \sim n/2$. A k -tuple of large dependent sets $B = (B_1, \dots, B_k)$ is simple, if for all sequences $(j_1 < j_2 < \dots < j_l)$ and $(1 \leq l \leq k)$ the set differences satisfy

$$|B_{j_1} \oplus B_{j_2} \oplus \dots \oplus B_{j_l}| \in J_1. \quad (17)$$

For any given matrix M there is a largest k such that B_1, \dots, B_k are simple. In which case, we say k is *maximal* and B_1, \dots, B_k is a *maximal simple sequence*.

Let $V(M) = \{\emptyset\} \cup \{B : B \text{ is zero-sum in } M\}$, then $(V(M), \oplus)$ is a vector space over GF_2 under the convention that $0 \cdot B = \emptyset$, $1 \cdot B = B$. In $V(M)$ a simple sequence (B_1, \dots, B_k) is an ordered basis for a subspace S of dimension k .

Given k linearly dependent sets of rows with index sets B_1, \dots, B_k , there are 2^k intersections of these sets and their complements. For each $\mathbf{x} = (x_1, \dots, x_k)$, $\mathbf{x} \in \{0, 1\}^k$ we let $I_{\mathbf{x}} = \bigcap_{i=1, \dots, k} B_i^{(x_i)}$ where $B_i^{(0)} = \overline{B}_i = [n] \setminus B_i$ and $B_i^{(1)} = B_i$. The index sets $I_{\mathbf{x}}$ are disjoint by definition and their union (including $\mathbf{x}_0 = (0, \dots, 0)$) is $[n]$.

Next for $\mathbf{x} \in \{0, 1\}^k$ let $B(\mathbf{x}) = \bigoplus_{i:x_i=1} B_i$. Let $K = 2^k - 1$. Let U be a $K \times K$ matrix indexed by $\mathbf{x}, \mathbf{y} \in \{0, 1\}^k$, $\mathbf{x}, \mathbf{y} \neq 0$; with entries $U(\mathbf{x}, \mathbf{y}) = 1$ if $I_{\mathbf{y}} \subseteq B(\mathbf{x})$, and $U(\mathbf{x}, \mathbf{y}) = 0$ otherwise. In summary,

$$\text{Row index } \mathbf{x} = (x_1, x_2, \dots, x_k) \text{ is the indicator vector for } B(\mathbf{x}) = \bigoplus_{i:x_i=1} B_i,$$

$$\text{Column index } \mathbf{y} = (y_1, y_2, \dots, y_k) \text{ is the indicator vector for } I_{\mathbf{y}} = \bigcap_{i=1, \dots, k} B_i^{(y_i)}.$$

The row of U representing the set $B(\mathbf{x})$ is formed by adding the rows of those sets B_i such that $x_i = 1$ in \mathbf{x} ; the addition being over $GF(2)$. Thus $B(\mathbf{x})$ is the union of the sets $I_{\mathbf{y}}$, where $y_i = 1$ for an odd number of those sets B_i where $x_i = 1$. This can be seen inductively by generating B_1 , $B_1 \oplus B_2$, $(B_1 \oplus B_2) \oplus B_3$ etc. in the given order. In summary $U(\mathbf{x}, \mathbf{y}) = 1$ iff both $x_i = 1$ and $y_i = 1$ for an odd number of indices i , and thus, over $GF(2)$,

$$U(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^k x_i y_i. \quad (18)$$

Our aim is to use U , treated as a real matrix to show that w.h.p. $|I_{\mathbf{x}}| \sim n/2^k$ for every \mathbf{x} . We do this by observing that given the characterisation $U(\mathbf{x}, \mathbf{y}) = 1_{I_{\mathbf{y}} \subseteq B(\mathbf{x})}$, the vector $(|I_{\mathbf{x}}|, \mathbf{x} \in \{0, 1\}^k, \mathbf{x} \neq 0)$ is the solution \mathbf{z} over the reals of an equation

$$U\mathbf{z} = \mathbf{b} \text{ where } \mathbf{b} \sim \frac{n}{2} \mathbf{1}, \quad (19)$$

assuming that $B = (B_1, \dots, B_k)$ is simple. To prove that $|I_{\mathbf{x}}| \sim n/2^k$, we prove the properties of U listed in Lemma 9 below.

Equation (18) implies that by arranging the rows and column indices of U in the same order, U will be symmetric. We will choose an ordering such the first k rows correspond to $B_i, i = 1, \dots, k$. Thus $x_i = e_i, i = 1, 2, \dots, k$ where $e_1 = (1, 0, \dots, 0)$ etc., and $y_i = e_i, i = 1, 2, \dots, k$. After this we let Q be the $k \times K$ matrix with column indices x made up of the first k rows. Thus row i represents $B_i, i = 1, \dots, k$ and U contains a $k \times k$ identity matrix in the first k rows and columns.

The row indexed by $\mathbf{x} = (x_1, \dots, x_k)$ is the linear combination $\sum_{i=1}^k x_i \mathbf{r}_i$ of the rows of Q , and corresponds to $B(\mathbf{x})$ in the vector space $V(M)$ given above.

Lemma 9. *The $K \times K$ matrix U has the following properties:*

- (i) The matrix U is symmetric.
- (ii) Every row or column of U has 2^{k-1} non-zero entries.
- (iii) Any two distinct rows of U have 2^{k-2} common non-zero entries.
- (iv) The matrix U is non-singular when the entries are taken to be over the real numbers, and the matrix $S = UU^\top = U^2 = 2^{k-2}(I + J)$ is symmetric, with inverse $S^{-1} = (1/2^{k-2})(I - J/2^k)$; where J is the all-ones matrix.

Proof. (i) This follows immediately from (18), and the above construction.

(ii) Fix \mathbf{x} and assume that $x_1 = 1$. There are 2^{k-1} choices for the values of $y_i, i = 2, 3, \dots, k$. Having made such a choice, there are two choices for y_1 , exactly one of which will give $\sum_{i=1}^k x_i y_i = 1$.

(iii) Fix \mathbf{x}, \mathbf{x}' and think of rows $\mathbf{x}, \mathbf{x}', \mathbf{x} + \mathbf{x}'$ as non-empty subsets of $[2^k]$. Then we have $|\mathbf{x}| = |\mathbf{x}'| = |\mathbf{x} \setminus \mathbf{x}'| + |\mathbf{x}' \setminus \mathbf{x}| = 2^{k-1}$, by (ii). Thus $|\mathbf{x}| + |\mathbf{x}'| - (|\mathbf{x} \setminus \mathbf{x}'| + |\mathbf{x}' \setminus \mathbf{x}|) = 2|\mathbf{x} \cap \mathbf{x}'| = 2^{k-1}$.

(iv) Let \mathbf{u}, \mathbf{v} be distinct rows of U , then $\mathbf{u} \cdot \mathbf{u} = 2^{k-1}$ and $\mathbf{u} \cdot \mathbf{v} = 2^{k-2}$. Thus $S = UU^\top = 2^{k-2}(I + J)$, where J is the all-ones matrix. By [3], Section 1.3, (1.9) and below, $\det(I + J) = 2^k \neq 0$, and thus S, U are non-singular. The reader can check that $S^{-1} = \frac{1}{2^{k-2}}(I - \frac{1}{2^k}J)$. \square

The definition of a simple k -tuple (B_1, \dots, B_k) requires that all sets B_i be large and their set differences to be distinct and of size $\sim n/2$. Let $(|B_1|, \dots, |B_k|) \sim (n/2)\mathbf{1}$ be the vector of these set sizes. Over the reals, solving (19) gives the sizes of the subsets $I_{\mathbf{x}}$.

Lemma 10. *Let (B_1, \dots, B_k) be a simple sequence. Then for all $\mathbf{x} \in \{0, 1\}^k$,*

$$|I_{\mathbf{x}}| = \frac{n}{2^k} \left(1 \pm 4^k \sqrt{\frac{\log n}{n}} \right). \quad (20)$$

Proof. Let $i = 1, \dots, K$ index the rows of U , and $j = 1, \dots, K$ index the columns. Let $B(i)$ be the set corresponding to the row i of U . Referring to (19), let $\mathbf{y} = (2/n)\mathbf{z}$, and $U\mathbf{y} = \mathbf{b}$ where now $b_i = 2|B(i)|/n = 1 + \varepsilon_i$, so that $|\varepsilon_i| \leq 2\sqrt{\log n/n}$. The matrix $S = U^2$, so $S\mathbf{y} = U\mathbf{b} = \mathbf{c}$ where $c_i = 2^{k-1}(1 + \delta_i)$ and $\delta_i = \sum_{j:U(i,j)=1} \varepsilon_j / 2^{k-1}$, the summation being over the 2^{k-1} -subset of non-zero entries of row i of U . Thus, as J is $K \times K$ where $K = 2^k - 1$,

$$\mathbf{y} = S^{-1}\mathbf{c} = \frac{1}{2^{k-2}} \left(I - \frac{1}{2^k}J \right) 2^{k-1}(\mathbf{1} + \boldsymbol{\delta}) = \frac{1}{2^{k-1}}\mathbf{1} + \boldsymbol{\eta},$$

where $|\boldsymbol{\eta}| \leq 2^k \sqrt{\log n/n}$. It follows that w.h.p. the solution \mathbf{z} to (19) over the real numbers satisfies $|I_{\mathbf{x}}| = (n/2^k)(1 \pm 4^k \sqrt{\log n/n})$ for all $\mathbf{x} \in \{0, 1\}^k$. \square

5 Simple sequences of large zero-sum sets

Let B_1, B_2, \dots, B_k be a simple sequence. In row M_i of the matrix M , there is a 1 in the diagonal entry $M_{i,i}$. As $s = 3$ there need to be two (random) 1's in column C_i chosen in a way to ensure the linear dependence of B_1, \dots, B_k . The following lemma describes where these non-zeros must be placed.

Lemma 11. B_1, \dots, B_k are dependencies if and only if the following holds for all $i \in [n]$. Suppose that row i is in $I_{\mathbf{x}}$, and that the two random non-zeros $e_1(i), e_2(i)$ in column i are in $I_{\mathbf{u}}, I_{\mathbf{v}}$ respectively. Then we must have $\mathbf{x} = \mathbf{u} + \mathbf{v} \pmod{2}$.

Proof. Let $\mathbf{x} = (x_1, \dots, x_k)$ and consider x_m for $1 \leq m \leq k$. If $x_m = 0$ then $i \notin B_m$, so either none or both of $e_1(i), e_2(i)$ are in B_m , and so zero or two unit entries in this column are in B_m . We must therefore have either $u_m = v_m = 0$ or $u_m = v_m = 1$ and $x_m = u_m + v_m$. If $x_m = 1$ then $i \in B_m$ and so exactly one of $e_1(i), e_2(i)$ must also be in B_m . Hence $u_m = 1, v_m = 0$, or vice versa. Thus in all cases $x_m = u_m + v_m$. \square

The main result of this section is the following.

Lemma 12. Let $k \geq 1$ be a positive integer, and let \mathbf{X}_k count the number of simple k -sequences of large dependencies. Then $\mathbf{E}(\mathbf{X}_k) \sim 1$.

Proof. We have to estimate the expected number of simple sequences (B_1, \dots, B_k) of large dependencies. By (20) of Lemma 10 the index sets $I_{\mathbf{x}}$ have size $|I_{\mathbf{x}}| = (n/2^k)(1 + O(\sqrt{\log n/n}))$. Let $K = 2^k - 1$ as above, and let

$$\Omega = \left\{ \mathbf{h} = (h_0, h_1, \dots, h_K) : h_i \text{ satisfies (20)}, \sum_{i=1}^K h_i \in J_1 \right\}.$$

Then we define $\Phi(\mathbf{h}, k)$ by

$$\mathbf{E}(\mathbf{X}_k) = \sum_{\mathbf{h} \in \Omega} \binom{n}{h_0, h_1, \dots, h_K} \prod_{\mathbf{x} \neq 0} \left(2 \sum_{\substack{\{\mathbf{u}, \mathbf{v}\} \\ \mathbf{u} + \mathbf{v} = \mathbf{x}}} \frac{h_{\mathbf{u}} h_{\mathbf{v}}}{n} \right)^{h_{\mathbf{x}}} \left(\sum_{\mathbf{u}} \left(\frac{h_{\mathbf{u}}}{n} \right)^2 \right)^{h_0} \quad (21)$$

$$= \sum_{\mathbf{h} \in \Omega} \binom{n}{h_0, h_1, \dots, h_K} \Phi(\mathbf{h}, k). \quad (22)$$

Explanation of (21). Let $h_{\mathbf{x}} = |I_{\mathbf{x}}|$. The multinomial coefficient $\binom{n}{h_0, h_1, \dots, h_K}$ counts the number of choices for the subsets $I_{\mathbf{x}}$. In the product, in order for B_1, \dots, B_k to be zero-sum, for $\mathbf{x} \neq 0$ we need to cancel the diagonal entries $M_{j,j} = 1$ of $\mathbf{j} \in I_{\mathbf{x}}$ within the columns

indexed by $I_{\mathbf{x}}$. This is achieved by putting one entry in rows $I_{\mathbf{u}}$ and one in rows $I_{\mathbf{v}}$ where $\mathbf{u} + \mathbf{v} = \mathbf{x}$. The last factor counts the choices for the entries of columns indexed by I_0 over the row index sets $I_{\mathbf{u}}$, either zero or two in an index set, in order to preserve the zero-sum property.

Set $h_{\mathbf{x}} = (n/2^k)(1 + \varepsilon_{\mathbf{x}})$ where $|\varepsilon_{\mathbf{x}}| = O(\sqrt{\log n/n})$. We note that $\sum_{\mathbf{x}} \varepsilon_{\mathbf{x}} = 0$, implies that

$$\sum_{\mathbf{x}} h_{\mathbf{x}} \varepsilon_{\mathbf{x}} = \frac{n}{2^k} \sum_{\mathbf{x}} (\varepsilon_{\mathbf{x}} + \varepsilon_{\mathbf{x}}^2) = \frac{n}{2^k} \sum_{\mathbf{x}} \varepsilon_{\mathbf{x}}^2 \text{ and } \sum_{\mathbf{x}} h_{\mathbf{x}} \varepsilon_{\mathbf{x}}^2 = \frac{n}{2^k} \sum_{\mathbf{x}} \varepsilon_{\mathbf{x}}^2 + O\left(\frac{\log^{3/2} n}{n^{1/2}}\right).$$

And then Stirling's approximation implies that

$$\begin{aligned} \binom{n}{h_0, h_1, \dots, h_K} &\sim \frac{n^n \sqrt{2\pi n}}{\prod_{\mathbf{x} \in \{0,1\}^k} ((n/2^k)(1 + \varepsilon_{\mathbf{x}}))^{h_{\mathbf{x}}} (\sqrt{2\pi n/2^k})^{2^k}} \\ &= 2^{kn} \exp \left\{ - \sum_{\mathbf{x} \in \{0,1\}^k} h_{\mathbf{x}} \left(\varepsilon_{\mathbf{x}} - \frac{\varepsilon_{\mathbf{x}}^2}{2} \right) + O(\log n) \right\} \\ &= 2^{kn} \exp \left\{ - \frac{n}{2^{k+1}} \sum_{\mathbf{x} \in \{0,1\}^k} \varepsilon_{\mathbf{x}}^2 + O(\log n) \right\} = 2^{kn} n^{O(1)}. \end{aligned}$$

In addition, by considering random 2^k -colorings of $[n]$ we see from the Chernoff bounds that

$$\sum_{\mathbf{h} \in \Omega} \binom{n}{h_0, h_1, \dots, h_K} = 2^{kn} (1 - O(n^{-2^k/3})). \quad (23)$$

With respect to (21), using $\sum_{\mathbf{x}} \varepsilon_{\mathbf{x}} = 0$, we see that

$$\begin{aligned} \left(\sum_{\mathbf{u} \in \{0,1\}^k} \left(\frac{h_{\mathbf{u}}}{n} \right)^2 \right)^{h_0} &= \left(\sum_{\mathbf{u}} \frac{1}{2^{2k}} (1 + 2\varepsilon_{\mathbf{u}} + \varepsilon_{\mathbf{u}}^2) \right)^{h_0} \\ &= \left(\frac{1}{2^k} \right)^{h_0} \left(1 + \frac{1}{2^k} \sum_{\mathbf{u}} \varepsilon_{\mathbf{u}}^2 \right)^{h_0} \\ &= \left(\frac{1}{2^k} \right)^{h_0} \exp \left\{ \frac{n}{2^k} (1 + \varepsilon_0) \log \left(1 + \sum_{\mathbf{u}} \frac{\varepsilon_{\mathbf{u}}^2}{2^k} \right) \right\} \\ &= \left(\frac{1}{2^k} \right)^{h_0} \exp \left\{ \frac{n}{2^{2k}} \sum_{\mathbf{u}} \varepsilon_{\mathbf{u}}^2 + O\left(\frac{\log^{3/2} n}{n^{1/2}}\right) \right\}. \end{aligned} \quad (24)$$

If $\mathbf{x} \neq 0$ then each index \mathbf{z} occurs exactly once in $\sum_{\{\mathbf{u}, \mathbf{v}\}} (\varepsilon_{\mathbf{u}} + \varepsilon_{\mathbf{v}})$ and so $\sum_{\{\mathbf{u}, \mathbf{v}\}} (\varepsilon_{\mathbf{u}} + \varepsilon_{\mathbf{v}}) = \sum_{\mathbf{z}} \varepsilon_{\mathbf{z}} = 0$. Therefore,

$$\begin{aligned}
\left(2 \sum_{\{\mathbf{u}, \mathbf{v}\}} \frac{h_{\mathbf{u}}}{n} \frac{h_{\mathbf{v}}}{n} \right)^{h_{\mathbf{x}}} &= \left(2 \sum_{\substack{\{\mathbf{u}, \mathbf{v}\} \\ \mathbf{u} + \mathbf{v} = \mathbf{x}}} \frac{1}{2^{2k}} (1 + \varepsilon_{\mathbf{u}} + \varepsilon_{\mathbf{v}} + \varepsilon_{\mathbf{u}} \varepsilon_{\mathbf{v}}) \right)^{h_{\mathbf{x}}} \\
&= \left(\frac{1}{2^k} \right)^{h_{\mathbf{x}}} \left(1 + \frac{1}{2^k} \sum_{\substack{\{\mathbf{u}, \mathbf{v}\} \\ \mathbf{u} + \mathbf{v} = \mathbf{x}}} 2\varepsilon_{\mathbf{u}} \varepsilon_{\mathbf{v}} \right)^{h_{\mathbf{x}}} \\
&= \left(\frac{1}{2^k} \right)^{h_{\mathbf{x}}} \exp \left\{ \frac{n}{2^k} (1 + \varepsilon_{\mathbf{x}}) \log \left(1 + 2 \sum_{\substack{\{\mathbf{u}, \mathbf{v}\} \\ \mathbf{u} + \mathbf{v} = \mathbf{x}}} \frac{\varepsilon_{\mathbf{u}} \varepsilon_{\mathbf{v}}}{2^k} \right) \right\} \\
&= \left(\frac{1}{2^k} \right)^{h_{\mathbf{x}}} \exp \left\{ \frac{n}{2^k} \sum_{\substack{\{\mathbf{u}, \mathbf{v}\} \\ \mathbf{u} + \mathbf{v} = \mathbf{x}}} \frac{2\varepsilon_{\mathbf{u}} \varepsilon_{\mathbf{v}}}{2^k} + O \left(\frac{\log^{3/2} n}{n^{1/2}} \right) \right\}.
\end{aligned}$$

Note that

$$\Lambda = \sum_{\mathbf{x} \neq 0} \sum_{\substack{\{\mathbf{u}, \mathbf{v}\} \\ \mathbf{u} + \mathbf{v} = \mathbf{x}}} 2\varepsilon_{\mathbf{u}} \varepsilon_{\mathbf{v}} = \sum_{\mathbf{u}} \varepsilon_{\mathbf{u}} \sum_{\substack{\mathbf{x} + \mathbf{u} \\ \mathbf{x} \neq 0}} \varepsilon_{\mathbf{x} + \mathbf{u}} = \sum_{\mathbf{u}} \varepsilon_{\mathbf{u}} \sum_{\mathbf{v} \neq \mathbf{u}} \varepsilon_{\mathbf{v}},$$

gives

$$\Lambda + \sum_{\mathbf{u}} \varepsilon_{\mathbf{u}}^2 = \left(\sum_{\mathbf{u}} \varepsilon_{\mathbf{u}} \right)^2 = 0.$$

Thus using $\sum_{\mathbf{x}} h_{\mathbf{x}} = n$,

$$\begin{aligned}
\Phi(\mathbf{h}, k) &= \left(\frac{1}{2^k} \right)^{\sum_{\mathbf{x}} h_{\mathbf{x}}} \exp \left\{ \frac{n}{2^{2k}} \left(\sum_{\mathbf{u}} \varepsilon_{\mathbf{u}}^2 + \sum_{\mathbf{x} \neq 0} \sum_{\substack{\{\mathbf{u}, \mathbf{v}\} \\ \mathbf{u} + \mathbf{v} = \mathbf{x}}} 2\varepsilon_{\mathbf{u}} \varepsilon_{\mathbf{v}} \right) + O \left(\frac{\log^{3/2} n}{n^{1/2}} \right) \right\} \\
&= \frac{1}{2^{kn}} e^{O(\log^{3/2} n / \sqrt{n})}.
\end{aligned} \tag{25}$$

It follows from (22), (23) and (25) above that

$$\mathbf{E}(\mathbf{X}_k) = 1 + O \left(\frac{\log^{3/2} n}{\sqrt{n}} \right) = 1 + o(1). \tag{26}$$

□

6 Conditional expected number of small zero-sum sets

Let (B_1, \dots, B_k) be a fixed sequence of subsets of $[n]$ with $|B_i| \in J_1$ for $i = 1, 2, \dots, k \leq \omega$. Let \mathcal{B} be the event

$$\mathcal{B} = \{(B_1, \dots, B_k) \text{ is a simple sequence of large row dependencies}\}. \quad (27)$$

Lemma 13. *Given \mathcal{B} and $i \in I_{\mathbf{x}}$, $|I_{\mathbf{x}}| = h_{\mathbf{x}}$, the distribution of the row indices ℓ, ℓ' of the other two non-zeros in column i is as follows.*

If $\mathbf{x} \neq 0$ then choose \mathbf{u}, \mathbf{v} such that $\mathbf{x} = \mathbf{u} + \mathbf{v} \pmod{2}$ with probability

$$p(\mathbf{u}, \mathbf{v}) = \frac{h_{\mathbf{u}} h_{\mathbf{v}}}{\sum_{\mathbf{y}+\mathbf{z}=\mathbf{x}} h_{\mathbf{y}} h_{\mathbf{z}}},$$

and then randomly choose $\ell \in I_{\mathbf{u}}, \ell' \in I_{\mathbf{v}}$. If $\mathbf{x} = 0$ then choose \mathbf{u} with probability

$$p(\mathbf{u}, \mathbf{u}) = \frac{h_{\mathbf{u}}^2}{\sum_{\mathbf{y} \in \{0,1\}^k} h_{\mathbf{y}}^2},$$

and then randomly choose $\ell, \ell' \in I_{\mathbf{u}}$.

Proof. This follows from the fact that the non-zeros in each column are independently chosen with replacement and from the condition given in Lemma 11. \square

For $m \leq \omega$, let S_j , $j = 1, 2, \dots, m$ be pairwise disjoint subsets of the rows of M , where $|S_j| \leq \omega$. Let $S = \bigcup_{j=1}^m S_j$ and $s = |S|$. For $j = 1, 2, \dots, m$ define the following events

$$\mathcal{S}_j = \{S_j \text{ is a small zero-sum set}\}, \quad \mathcal{S}_j^* = \{S_j \text{ is a small fundamental zero-sum set}\}.$$

Let

$$\mathcal{S} = \bigcap_{j=1}^m \mathcal{S}_j \quad \text{and} \quad \mathcal{S}^* = \bigcap_{j=1}^m \mathcal{S}_j^*.$$

We need to understand the conditioning imposed by the event \mathcal{B} in (27) on the small dependencies.

Lemma 14.

$$\mathbb{P}(\mathcal{S}^* \mid \mathcal{B}) \sim \mathbb{P}(\mathcal{S}^*). \quad (28)$$

Proof. Let $I_{\mathbf{x}}$, $\mathbf{x} \in \{0,1\}^k$, be as defined in Section 4. Let $h_{\mathbf{x}} = |I_{\mathbf{x}}|$. By Lemma 10 we can assume that $|I_{\mathbf{x}}| = h_{\mathbf{x}} \sim n/2^k$ for all $\mathbf{x} \in \{0,1\}^k$. For $j = 1, 2, \dots, m$, let $S_{j,\mathbf{x}} = S_j \cap I_{\mathbf{x}}$ and $s_{j,\mathbf{x}} = |S_{j,\mathbf{x}}|$. Similarly, let $S_{\mathbf{x}} = S \cap I_{\mathbf{x}}$, $s_{\mathbf{x}} = |S_{\mathbf{x}}|$. These definitions include $\mathbf{x} = \mathbf{0}$, so that $S_{\mathbf{0}} = I_{\mathbf{0}} \cap S$ and $s_{j,\mathbf{0}} = |S_{j,\mathbf{0}}|$ etc.

For each $i \in [n]$, we consider the probability that column i of M is consistent with \mathcal{S} according to four cases.

Case 1: $i \in I_0 \setminus S$. For each column $i \in I_0 \setminus S = I_0 \setminus S_0$, we must estimate the probability that the two non-zeros $e_1(i), e_2(i)$ are in rows consistent with the occurrence of \mathcal{S} . Because $i \in I_0$ and \mathcal{B} occurs, we know from Lemma 11 that $e_1(i), e_2(i) \in I_u$ for some $\mathbf{u} \in \{0, 1\}^k$. For \mathcal{S} to occur, we require that zero or two of $e_1(i), e_2(i)$ fall in S_u , an event of conditional probability $(1 - s_u/h_u)^2 + (s_u/h_u)^2$.

Let E_u denote the number of non-zero pairs from $I_0 \setminus S_0$ falling in I_u . Then the conditional probability that the non-zeros of $I_0 \setminus S_0$ are consistent with \mathcal{S} is given by

$$\mathbb{P}(I_0 \setminus S_0 \text{ is consistent with } \mathcal{S} \mid \mathcal{B}) = \mathbb{E} \left(\prod_u \left(1 - 2 \frac{s_u}{h_u} + 2 \left(\frac{s_u}{h_u} \right)^2 \right)^{E_u} \right). \quad (29)$$

Given \mathcal{B} , we see that E_u is distributed as $\text{Bin}(h_0 - s_0, p(\mathbf{u}, \mathbf{u}))$, and has expectation

$$\mathbb{E}(E_u) = (h_0 - s_0) \frac{h_u^2}{h_0^2 + h_1^2 + \dots + (h_{2^k-1})^2} \sim \frac{h_0}{2^k}.$$

By Lemma 10 we can assume that $h_0 \sim N = n/2^k$. The Chernoff bounds imply that E_u is concentrated around its mean $(h_0 - s_0)p(\mathbf{u}, \mathbf{u})$. Thus,

$$\left| E_u - \frac{h_0}{2^k} \right| \leq n^{2/3} \quad \text{with probability at least } 1 - e^{-\Omega(n^{1/3})}. \quad (30)$$

Going back to (29) and using (30) gives

$$\begin{aligned} \mathbb{P}(I_0 \setminus S_0 \text{ is consistent with the occurrence of } \mathcal{S} \mid \mathcal{B}) &\sim \\ \prod_u \left(1 - \frac{2s_u}{N} \right)^{N/2^k} &\sim \exp \left\{ -2 \sum_u \frac{s_u}{2^k} \right\} = e^{-s/2^{k-1}}. \end{aligned} \quad (31)$$

Case 2: $i \in I_x \setminus S$, $x \neq 0$. Given \mathcal{B} , and $i \in I_x$, we know from Lemma 11 that the non-zeros $e_1(i), e_2(i)$ of column i lie in I_u, I_{x+u} respectively, for some $\mathbf{u} \in \{0, 1\}^k$. The probability of this is $p(\mathbf{u}, \mathbf{x} + \mathbf{u})$. The number $E_x(\mathbf{u}, \mathbf{x} + \mathbf{u})$ of such pairs of non-zeros in I_u, I_{x+u} has distribution $\text{Bin}((h_x - s_x)p(\mathbf{u}, \mathbf{x} + \mathbf{u}))$, and expectation asymptotic to $(h_x - s_x)/2^{k-1}$.

The rows of S_1, \dots, S_m have to be zero-sum in this column, so either exactly one non-zero falls in **some** $S_{j,u}, S_{j,x+u}$ for some $1 \leq j \leq m$ or exactly one non-zero falls in **some**

$I_{\mathbf{u}} \setminus S_{\mathbf{u}}, I_{\mathbf{x}+\mathbf{u}} \setminus S_{\mathbf{x}+\mathbf{u}}$. The **conditional** probability of this is

$$\begin{aligned} P(\mathbf{u}, \mathbf{x} + \mathbf{u}) &= \mathbf{E} \left(\left(\sum_{j=1}^m \frac{s_{j,\mathbf{u}}}{h_{\mathbf{u}}} \frac{s_{j,\mathbf{x}+\mathbf{u}}}{h_{\mathbf{x}+\mathbf{u}}} + \frac{h_{\mathbf{u}} - s_{\mathbf{u}}}{h_{\mathbf{u}}} \frac{h_{\mathbf{x}+\mathbf{u}} - s_{\mathbf{x}+\mathbf{u}}}{h_{\mathbf{x}+\mathbf{u}}} \right)^{E_{\mathbf{x}}(\mathbf{u}, \mathbf{x} + \mathbf{u})} \right) \\ &\sim \left(\sum_{j=1}^m \frac{s_{j,\mathbf{u}} s_{j,\mathbf{x}+\mathbf{u}}}{N^2} + \frac{N - s_{\mathbf{u}}}{N} \frac{N - s_{\mathbf{x}+\mathbf{u}}}{N} \right)^{(N - s_{\mathbf{x}})/2^{k-1}} \\ &\sim e^{-(s_{\mathbf{u}} + s_{\mathbf{x}+\mathbf{u}})/2^{k-1}}. \end{aligned}$$

For a given \mathbf{x} there are 2^{k-1} unordered pairs $S_{\mathbf{u}}, S_{\mathbf{x}+\mathbf{u}}$, so

$$\mathbb{P}(I_{\mathbf{x}} \setminus S_{\mathbf{x}} \text{ is consistent with } \mathcal{S}) \sim \exp \left\{ -\frac{1}{2^{k-1}} \sum_{\{\mathbf{u}, \mathbf{x}+\mathbf{u}\}} (s_{\mathbf{u}} + s_{\mathbf{x}+\mathbf{u}}) \right\} = e^{-s/2^{k-1}}. \quad (32)$$

Note that, in the sum in (32) $s_{\mathbf{u}} + s_{\mathbf{x}+\mathbf{u}}$ and $s_{\mathbf{x}+\mathbf{u}} + s_{\mathbf{u}}$, contribute as one term. Thus

$$\mathbb{P}(I_{\mathbf{x}} \setminus S_{\mathbf{x}} \text{ is consistent with } \mathcal{S}, \forall \mathbf{x} \neq \mathbf{0}) \sim e^{-(2^{k-1}-1)s/2^{k-1}}. \quad (33)$$

Case 3: $i \in S_{j,\mathbf{x}} \subseteq I_{\mathbf{x}}$, $\mathbf{x} \neq 0$. Suppose that the pair $e_1(i), e_2(i)$ fall in $I_{\mathbf{u}}, I_{\mathbf{u}+\mathbf{x}}$. For $i \in S_{j,\mathbf{x}}$, one non-zero needs to be in S_j , and the other to **completely avoid** \mathcal{S} . Let $\mathbf{v} = \mathbf{x} + \mathbf{u}$. The probability this happens is

$$P_j(\mathbf{u}, \mathbf{v}) \sim \frac{1}{2^{k-1}} \left(\frac{s_{j,\mathbf{u}}}{h_{\mathbf{u}}} \frac{h_{\mathbf{v}} - s_{\mathbf{v}}}{h_{\mathbf{v}}} + \frac{s_{j,\mathbf{v}}}{h_{\mathbf{v}}} \frac{h_{\mathbf{u}} - s_{\mathbf{u}}}{h_{\mathbf{u}}} \right). \quad (34)$$

The events $\{\mathbf{u}, \mathbf{x} + \mathbf{u}\}$ are disjoint and **are an exhaustive dissection of** S_j . For a given $i \in S_{j,\mathbf{x}}$, the probability $p(i, j)$ of success is

$$\begin{aligned} p(i, j) &= \sum_{\{\mathbf{u}, \mathbf{u}+\mathbf{x}\}} P_j(\mathbf{u}, \mathbf{u}+\mathbf{x}) \sim \frac{1}{2^{k-1}} \sum_{\mathbf{u}, \mathbf{v}=\mathbf{x}+\mathbf{u}} \left(\frac{s_{j,\mathbf{u}}}{N} \frac{N - s_{\mathbf{v}}}{N} + \frac{s_{j,\mathbf{v}}}{N} \frac{N - s_{\mathbf{u}}}{N} \right) \\ &\sim \frac{s_j}{N 2^{k-1}} \left(1 + O\left(\frac{\omega}{N}\right) \right). \end{aligned} \quad (35)$$

Every column of $S_{j,\mathbf{x}}$ has to succeed or **some** S_t is not a small zero-sum set. Thus

$$\mathbb{P}(S_{j,\mathbf{x}} \text{ succeeds}) \sim \left(\frac{s_j (1 + O(s/N))}{N 2^{k-1}} \right)^{s_{j,\mathbf{x}}}.$$

As $\sum_{\mathbf{x} \neq \mathbf{0}} s_{j,\mathbf{x}} = s_j - s_{j,\mathbf{0}}$, the above allows us to calculate

$$\mathbb{P}(S_{j,\mathbf{x}} \text{ succeeds } \forall \mathbf{x} \neq \mathbf{0}) \sim \left(\frac{s_j}{N 2^{k-1}} \right)^{s_j - s_{j,\mathbf{0}}}. \quad (36)$$

Case 4: $i \in S_{j,0} \subseteq \mathcal{I}_0$. In the case that $\mathbf{x} = \mathbf{0}$, and $S_{j,0} \subseteq \mathcal{I}_0$, the non-zeros in a column of $S_{j,0}$ must both fall in the same index set $I_{\mathbf{u}}$; one in $S_{j,\mathbf{u}}$ and one in $I_{\mathbf{u}} \setminus S_{j,\mathbf{u}}$. Thus $P(\mathbf{u}, \mathbf{u})$ is now summed over all $I_{\mathbf{u}}$, a total of 2^k such sets. For $i \in S_{j,0}$, the probability $p(i)$ of success is

$$p(i) = \sum_{\{\mathbf{u}, \mathbf{u}\}} P(\mathbf{u}, \mathbf{u}) \sim \frac{1}{2^k} \sum_{\mathbf{u}} \left(2 \frac{s_{j,\mathbf{u}}}{N} \frac{N - s_{j,\mathbf{u}}}{N} \right) \sim \frac{s_j}{N 2^{k-1}} \left(1 + O\left(\frac{\omega}{N}\right) \right).$$

The final term is the same as in (35), and we obtain

$$\mathbb{P}(S_{j,0} \text{ succeeds}) \sim \left(\frac{s_j}{N 2^{k-1}} \right)^{s_{j,0}} \quad (37)$$

Using (31), (33), (36) and (37), we obtain

$$\mathbb{P}(\mathcal{S} \mid \mathcal{B}) \sim \prod_{j=1}^m \left(\frac{s_j}{N 2^{k-1}} \right)^{s_j} e^{-(2^{k-1})s/2^{k-1}} e^{-s/2^{k-1}} = \prod_{j=1}^m \left(\frac{2s_j}{n} \right)^{s_j} e^{-2s}. \quad (38)$$

Applying (6) to the right hand side of (38) completes the proof of $\mathbb{P}(\mathcal{S} \mid \mathcal{B}) \sim \mathbb{P}(\mathcal{S})$. To replace \mathcal{S} by \mathcal{S}^* the conditional probability that S_j is fundamental is obtained by multiplying by κ_{s_j} of (7). This completes the proof of the lemma. \square

We can now use inclusion-exclusion to prove the following lemma.

Lemma 15. *Let Σ_{σ} be the event that there are exactly σ disjoint small fundamental dependencies. Then,*

$$\mathbb{P}(\Sigma_{\sigma} \mid \mathcal{B}) \sim \frac{\phi_R^{\sigma} e^{-\phi_R}}{\sigma!} \sim \mathbb{P}(\Sigma_{\sigma}).$$

Proof. Let $s = s_1 + \dots + s_{\ell}$, then

$$\begin{aligned} T_{\ell} &= \frac{1}{\ell!} \sum_{1 \leq s_1, \dots, s_{\ell} \leq \omega} \sum_{\substack{|S_i|=s_i, \\ i=1, \dots, \ell}} \mathbb{P} \left(\bigcap_{i=1}^{\ell} \mathcal{S}_i^* \mid \mathcal{B} \right) \sim \frac{1}{\ell!} \sum_{1 \leq s_1, \dots, s_{\ell} \leq \omega} \sum_{\substack{|S_i|=s_i, \\ i=1, \dots, \ell}} \mathbb{P} \left(\bigcap_{i=1}^{\ell} \mathcal{S}_i^* \right) \\ &\sim \frac{1}{\ell!} \sum_{1 \leq s_1, \dots, s_{\ell} \leq \omega} \binom{n}{s_1, \dots, s_{\ell}, n-s} \prod_{i=1}^{\ell} \left(\frac{2s_i}{n} \right)^{s_i} e^{-2s_i} \kappa_{s_i} \sim \frac{1}{\ell!} \sum_{1 \leq s_1, \dots, s_{\ell} \leq \omega} \prod_{i=1}^{\ell} \frac{(2s_i)^{s_i}}{s_i!} e^{-2s_i} \kappa_{s_i} \\ &\sim \frac{1}{\ell!} \left(\sum_{s=1}^{\infty} \frac{(2e^{-2})^s}{s} \sigma_s \right)^{\ell} \sim \frac{\phi_R^{\ell}}{\ell!}. \end{aligned}$$

The first approximation follows from Lemma 14 and the second from (6), (7).

Using Inclusion-Exclusion, we have

$$\mathbb{P}(\Sigma_\sigma \mid \mathcal{B}) = \sum_{\ell \geq \sigma} (-1)^{\ell-\sigma} \binom{\ell}{\sigma} T_\ell \sim \sum_{\ell \geq \sigma} (-1)^{\ell-\sigma} \binom{\ell}{\sigma} \frac{\phi_R^\ell}{\ell!} = \frac{\phi_R^\sigma e^{-\phi_R}}{\sigma!}.$$

Lemma 7 gives the unconditional probability. \square

Let \mathbf{X}_k count the number of simple k -sequences as in Lemma 12.

Lemma 16. *If $\sigma = O(1)$ then $\mathbf{E}(\mathbf{X}_k \mid \Sigma_\sigma) \sim 1$.*

Proof.

$$\begin{aligned} \mathbf{E}(\mathbf{X}_k \mid \Sigma_\sigma) &= \sum_{\mathcal{B}=(B_1, \dots, B_k)} \mathbb{P}(\mathcal{B} \mid \Sigma_\sigma) \\ &= \sum_{\mathcal{B}=(B_1, \dots, B_k)} \frac{\mathbb{P}(\Sigma_\sigma \mid \mathcal{B}) \mathbb{P}(\mathcal{B})}{\mathbb{P}(\Sigma_\sigma)} \\ &= \sum_{\mathcal{B}=(B_1, \dots, B_k)} \frac{\mathbb{P}(\mathcal{B})}{\mathbb{P}(\Sigma_\sigma)} \sum_{\ell \geq \sigma} (-1)^{\ell-\sigma} \binom{\ell}{\sigma} T_\ell \\ &= \sum_{\mathcal{B}=(B_1, \dots, B_k)} \frac{\mathbb{P}(\mathcal{B})}{\mathbb{P}(\Sigma_\sigma)} \sum_{\ell \geq \sigma} (-1)^{\ell-\sigma} \binom{\ell}{\sigma} \frac{1}{\ell!} \sum_{1 \leq s_1, \dots, s_\ell \leq \omega} \sum_{\substack{|S_i|=s_i, \\ i=1, \dots, \ell}} \mathbb{P}\left(\bigcap_{i=1}^{\ell} S_i^* \mid \mathcal{B}\right) \\ &\sim \sum_{\mathcal{B}=(B_1, \dots, B_k)} \frac{\mathbb{P}(\mathcal{B})}{\mathbb{P}(\Sigma_\sigma)} \sum_{\ell \geq \sigma} (-1)^{\ell-\sigma} \binom{\ell}{\sigma} \frac{1}{\ell!} \sum_{1 \leq s_1, \dots, s_\ell \leq \omega} \sum_{\substack{|S_i|=s_i, \\ i=1, \dots, \ell}} \mathbb{P}\left(\bigcap_{i=1}^{\ell} S_i^*\right) \\ &\sim \sum_{\mathcal{B}=(B_1, \dots, B_k)} \frac{\mathbb{P}(\mathcal{B})}{\mathbb{P}(\Sigma_\sigma)} \mathbb{P}(\Sigma_\sigma) \\ &= \mathbf{E}(\mathbf{X}_k) \sim 1. \end{aligned}$$

\square

7 Joint distribution of small and large dependencies

We first state a preparatory lemma. A proof of the next result for $c_k = 1$ can be found in [6], [7]. We give a full and different proof for completeness.

Lemma 17. For $\lambda \geq 0$, and $k \geq 0$ the solutions to

$$c_k = \sum_{\lambda=k}^{\infty} q_{\lambda} \prod_{i=0}^{k-1} (2^{\lambda} - 2^i), \quad (39)$$

are given by

$$q_{\lambda} = \sum_{k=\lambda}^{\infty} (-1)^{k-\lambda} 2^{\binom{k-\lambda}{2}} \begin{bmatrix} k \\ \lambda \end{bmatrix}_2 \psi_k c_k, \quad (40)$$

where $\psi_k = 1 / \left(2^{\binom{k}{2}} \prod_{i=1}^k (2^i - 1) \right)$. In particular, if $c_k = 1$, $q_{\lambda} = \pi(\lambda)$ of (4).

Proof. Gaussian coefficients are defined as

$$\begin{bmatrix} \lambda \\ k \end{bmatrix}_z = \frac{\prod_{i=1}^k (z^{\lambda-i+1} - 1)}{\prod_{i=1}^k (z^i - 1)}. \quad (41)$$

Using (41) with $z = 2$, equation (39) can be rewritten as

$$c_k = 2^{\binom{k}{2}} \prod_{i=1}^k (2^i - 1) \sum_{\lambda=k}^{\infty} q_{\lambda} \begin{bmatrix} \lambda \\ k \end{bmatrix}_2. \quad (42)$$

Put $\psi_k = 1 / \left(2^{\binom{k}{2}} \prod_{i=1}^k (2^i - 1) \right)$, we see that q_{λ} is the solution to

$$\sum_{\lambda=k}^{\infty} \begin{bmatrix} \lambda \\ k \end{bmatrix}_2 q_{\lambda} = \psi_k c_k, \quad k \geq 0. \quad (43)$$

Fix $\delta \geq 0$, multiply equation $k \geq \delta$ in (43) by $(-1)^{k-\delta} 2^{\binom{k-\delta}{2}} \begin{bmatrix} k \\ \delta \end{bmatrix}_2$, and sum these equations over $k \geq \delta$. This gives

$$\sum_{k=\delta}^{\infty} (-1)^{k-\delta} 2^{\binom{k-\delta}{2}} \begin{bmatrix} k \\ \delta \end{bmatrix}_2 \psi_k c_k = \sum_{k=\delta}^{\infty} \sum_{\lambda=k}^{\infty} (-1)^{k-\delta} \begin{bmatrix} k \\ \delta \end{bmatrix}_2 2^{\binom{k-\delta}{2}} \begin{bmatrix} \lambda \\ k \end{bmatrix}_2 q_{\lambda} \quad (44)$$

$$= \sum_{k=\delta}^{\infty} \sum_{\lambda=k}^{\infty} (-1)^{k-\delta} \begin{bmatrix} \lambda - \delta \\ k - \delta \end{bmatrix}_2 2^{\binom{k-\delta}{2}} \begin{bmatrix} \lambda \\ \delta \end{bmatrix}_2 q_{\lambda} \quad (45)$$

$$= \sum_{\lambda=\delta}^{\infty} \begin{bmatrix} \lambda \\ \delta \end{bmatrix}_2 q_{\lambda} \sum_{k=\delta}^{\lambda} (-1)^{k-\delta} \begin{bmatrix} \lambda - \delta \\ k - \delta \end{bmatrix}_2 2^{\binom{k-\delta}{2}} \quad (46)$$

Explanation: (45) to (46): Gaussian coefficients satisfy the identity

$$(1+x)(1+zx)\cdots(1+z^{r-1}x) = \sum_{\ell=0}^r \begin{bmatrix} r \\ \ell \end{bmatrix}_z z^{(\ell)} x^\ell. \quad (47)$$

To prove the last summation on the right hand side of (45) is zero for $\lambda > \delta$, use (47) with $x = -1, z = 2, \ell = k - \delta$ and $r = \lambda - \delta$. This gives $\sum_{\ell=0}^{\lambda-\delta} \begin{bmatrix} \lambda-\delta \\ \ell \end{bmatrix}_2 2^{(\ell)} (-1)^\ell = 0$ for $\lambda > \delta$.

For $z < 1$, taking the limit of (47) gives

$$\prod_{\ell=0}^{\infty} (1+z^\ell x) = \sum_{\ell=0}^{\infty} \frac{z^{(\ell)} x^\ell}{\prod_{i=1}^{\ell} (1-z^i)}. \quad (48)$$

Replacing δ by λ , and putting $c_k = 1$ in (40), we see that the solution q_λ to (39) is

$$\begin{aligned} q_\lambda &= \sum_{k=\lambda}^{\infty} \frac{(-1)^{k-\lambda} 2^{\binom{k-\lambda}{2} - \binom{k}{2}}}{\prod_{i=0}^{\lambda-1} (2^{\lambda-i} - 1) \prod_{i=\lambda}^{k-1} (2^{k-i} - 1)} \\ &= \frac{\left(\frac{1}{2}\right)^{\lambda^2}}{\prod_{i=1}^{\lambda} \left(1 - \left(\frac{1}{2}\right)^i\right)} \sum_{\ell=0}^{\infty} \frac{(-1)^\ell \left(\frac{1}{2}\right)^{\binom{\ell}{2}} \left(\frac{1}{2}\right)^{(1+\lambda)\ell}}{\prod_{i=1}^{\ell} \left(1 - \left(\frac{1}{2}\right)^i\right)} \end{aligned} \quad (49)$$

$$= \left(\frac{1}{2}\right)^{\lambda^2} \frac{\prod_{i=\lambda+1}^{\infty} \left(1 - \left(\frac{1}{2}\right)^i\right)}{\prod_{i=1}^{\lambda} \left(1 - \left(\frac{1}{2}\right)^i\right)} = \pi(\lambda), \quad (50)$$

where $\pi(\lambda)$ is given in (4). To get from (49) to (50), use (48) with $z = 1/2$ and $x = (-1/2^{\lambda+1})$. \square

Quotient space argument

Given M , let $\mathcal{B} = \{B_i : i \in [N]\}$ denote the set of large dependencies and $\mathcal{S} = \{S_j : j \in [T]\}$ denote the set of small dependencies. The following observations complete the proof of Theorem 1.

P1 Suppose that V, V_S are the vector spaces generated by all dependencies, and small dependencies, respectively. Suppose that these spaces have dimensions d, σ respectively.

Let $W = V/V_S$ be the quotient space and f_S be the canonical map $f_S : V \rightarrow W$. Thus f_S maps small dependencies to zero and $W = \{f_S(B) : B \in \mathcal{B}\} \cup \{0\}$. Each vector in W corresponds to an equivalence class of vectors in V . In terms of dependencies in \mathcal{B} , $B \sim B'$ iff $B \oplus B' = S$ where $S \in \mathcal{S}$. As the small dependencies are disjoint, the size of the equivalence class of B is 2^σ .

P2 Note that $\dim(W) = \dim(V) - \dim(V_S) = d - \sigma$. Let λ denote the maximum number of independent large dependencies. This will be the same as the maximum length of a simple sequence. We next prove that $\lambda = \dim(W)$.

Let $\mathbf{b}_i, i = 1, 2, \dots, m$ be a basis of W then $B_i \in f_S^{-1}(\mathbf{b}_i), i = 1, 2, \dots, m$ form a simple sequence. If not then for some $A \subseteq [m]$ we have $\bigoplus_{i \in A} B_i \in V_S$ which implies that $f_S(\bigoplus_{i \in A} B_i) = \sum_{i \in A} \mathbf{b}_i = 0$. Conversely, if B_1, B_2, \dots, B_k is a simple sequence then $\mathbf{b}_i = f_S(B_i), i = 1, 2, \dots, k$ are independent. If not then for some $A \subseteq [k]$, $\sum_{i \in A} \mathbf{b}_i = 0$ which implies that $\bigoplus_{i \in A} B_i \in V_S$.

P3 The first i independent members of a simple sequence generate a vector space W_i of size 2^i . The next independent entry of the sequence is chosen from $W \setminus W_i$, a space of size $2^\lambda - 2^i$. Each entry is chosen from an equivalence class of size 2^σ . It follows that the number X_k of simple sequences of length k is equal to

$$\prod_{i=0}^{k-1} ((2^\lambda - 2^i) \times 2^\sigma) = 2^{k\sigma} \prod_{i=0}^{k-1} (2^\lambda - 2^i).$$

P4 Let $b_t = \mathbb{P}(\lambda = t \mid \sigma = s)$. By Lemma 16, $\mathbf{E}(X_k \mid \sigma = s) \sim 1$, so

$$1 \sim \mathbf{E}(X_k \mid \sigma = s) = 2^{sk} \sum_{t=k}^{\infty} \prod_{i=0}^{k-1} (2^t - 2^i) b_t. \quad (51)$$

This can be re-written (with \sim replaced by $=$) as,

$$2^{-sk} = 2^{\binom{k}{2}} \prod_{i=1}^k (2^i - 1) \sum_{t=k}^{\infty} b_t \begin{bmatrix} t \\ k \end{bmatrix}_2$$

By Lemma 17 we find that

$$\begin{aligned}
b_t &= \sum_{k=t}^{\infty} (-1)^{k-t} 2^{\binom{k-t}{2}} \begin{bmatrix} k \\ t \end{bmatrix}_2 \psi_k c_k \\
&= \sum_{k=t}^{\infty} \frac{(-1)^{k-t} 2^{\binom{k-t}{2} - \binom{k}{2} - ks}}{\prod_{i=1}^k (2^i - 1)} \begin{bmatrix} k \\ t \end{bmatrix}_2 \\
&= \frac{1}{(2^t - 1) \cdots (2 - 1)} \sum_{k \geq t} (-1)^{k-t} 2^{\binom{k-t}{2} - \binom{k}{2} - ks - \binom{k+1-t}{2}} \frac{1}{\prod_{i=1}^{k-t} (1 - (1/2)^i)} \\
&= \left(\frac{1}{2}\right)^{t(t+s)} \frac{1}{\prod_{j=1}^t (1 - 1/2^j)} \sum_{j \geq 0} \left(\frac{1}{2}\right)^{\binom{j}{2}} \left(-1 \left(\frac{1}{2}\right)^{1+s+t}\right)^j \frac{1}{\prod_{i=1}^j (1 - (1/2)^i)} \\
&= \left(\frac{1}{2}\right)^{t(t+s)} \frac{1}{\prod_{j=1}^t (1 - 1/2^j)} \prod_{j=0}^{\infty} \left(1 - \left(\frac{1}{2}\right)^{(s+t+1)+j}\right) \\
&= P(s, t),
\end{aligned}$$

as given in (2), and where we used (48) with $z = 1/2$ and $x = -(1/2)^{s+t+1}$ to replace the alternating sum.

P5 The $P(s, t)$ only satisfy the solution $b_t(s) = \mathbb{P}(\lambda = t \mid \sigma = s)$ in (51) asymptotically. So to prove the lemma, we show that for large K ,

$$\sum_{\substack{t \geq K \\ s \geq 0}} b_t(s) \leq \varepsilon, \quad (52)$$

where $\varepsilon > 0$ is arbitrarily small. For $t \geq k$,

$$\prod_{i=0}^{k-1} (2^t - 2^i) = 2^{kt} \prod_{i=0}^{k-1} \left(1 - \frac{1}{2^{t-i}}\right) \geq 2^{kt} \left(1 - \sum_{i=0}^{k-1} \frac{1}{2^{t-i}}\right) \geq 2^{(k-1)t}.$$

It follows that

$$\sum_{\substack{t \geq K \\ s \geq 0}} b_t(s) \leq 2^{-K(K-1)}.$$

Thus (52) holds if $K \geq \sqrt{2 \log_2 1/\varepsilon}$.

References

[1] D. Achlioptas and M. Molloy. The solution space geometry of random linear equations, *Random Structures and Algorithms* 46.2, 197–231, (2015).

- [2] B. Bollobas. *Random Graphs*, 2nd edition. Cambridge University Press (2001).
- [3] R. Brualdi and H. Ryser. *Combinatorial Matrix Theory*. Cambridge University Press. (1991).
- [4] T. Bohman and A.M. Frieze. Hamilton cycles in 3-out, *Random Structures and Algorithms* 35, 393-417, (2009).
- [5] A. Coja-Oghlan, A. Ergür, P. Gao, S. Hetterich, M. Rolvien. The rank of sparse random matrices, SODA 2020, 579-591, (2020).
- [6] C. Cooper. On the rank of random matrices, *Random Structures and Algorithms* 16, 209-232, (2000).
- [7] C. Cooper. On the distribution of rank of a random matrix over a finite field, *Random Structures and Algorithms* 17, 197-212, (2000).
- [8] C. Cooper, A.M. Frieze and W. Pegden. On the rank of a random binary matrix, SODA 2019, 946-955, (2019).
- [9] C. Cooper and A.M. Frieze. Rank of the vertex-edge incidence matrix of r -out hypergraphs. Extended ArXiv version (2021). <https://arxiv.org/pdf/2107.05779.pdf>
- [10] T. Fenner and A.M. Frieze. On the connectivity of random m-orientable graphs and digraphs, *Combinatorica* 2, 347-359, (1982).
- [11] A.M. Frieze. Maximum matchings in a class of random graphs, *Journal of Combinatorial Theory B* 40, 196-212, (1986).
- [12] A.M. Frieze and M. Karoński. *Introduction to Random Graphs*, Cambridge University Press, (2016).
- [13] M. Ibrahimi, Y. Kanoria, M. Kraning and A. Montanari. The set of solutions of random XORSAT formulae, *Annals of Applied Probability*, 25.5, 2743–2808, (2015).
- [14] I. N. Kovalenko, A. A. Levitskya and M. N. Savchuk. *Selected Problems in Probabilistic Combinatorics*. Naukova Dumka, Kyiv (1986) (in Russian).