

Intercalates and Discrepancy in Random Latin Squares

Matthew Kwan *

Benny Sudakov[†]

Abstract

An *intercalate* in a Latin square is a 2×2 Latin subsquare. Let \mathbf{N} be the number of intercalates in a uniformly random $n \times n$ Latin square. We prove that asymptotically almost surely $\mathbf{N} \geq (1 - o(1)) n^2/4$, and that $\mathbb{E}\mathbf{N} \leq (1 + o(1)) n^2/2$ (therefore asymptotically almost surely $\mathbf{N} \leq fn^2$ for any $f \rightarrow \infty$). This significantly improves the previous best lower and upper bounds. We also give an upper tail bound for the number of intercalates in two fixed rows of a random Latin square.

In addition, we give a short proof of an asymptotically almost sure upper bound on the discrepancy of a random Latin square. This constitutes progress towards a conjecture of Linial and Luria.

1 Introduction

A $n \times n$ *Latin square* is an $n \times n$ array of the numbers between 1 and n (we call these *symbols*), such that each row and column contains each symbol exactly once. Latin squares are a fundamental type of combinatorial design, and have many essentially equivalent formulations (see for example [3, Section III, Theorem 1.11]).

The uniform probability distribution over the set \mathcal{L} of $n \times n$ Latin squares is poorly understood, due to the rigidity of a Latin square and the lack of any independence or recursive structure. It is not even known how to efficiently generate a uniformly random $\mathbf{L} \in \mathcal{L}$. Jacobson and Matthews [5] and Pittenger [13] designed Markov chains on \mathcal{L} which converge to the uniform distribution, but it is not known if these Markov chains converge rapidly. One obstacle to empirically testing the speed of convergence is that little is known about the likely structure of a random Latin square.

An *intercalate* in a Latin square L is a 2×2 Latin subsquare. That is, it is a pair of rows i, j and a pair of columns x, y such that $L_{i,x} = L_{j,y}$ and $L_{i,y} = L_{j,x}$. One important statistic of a Latin square L is the number $N(L)$ of intercalates that it contains. Clearly this number is at most n^3 , because each of the n^2 entries in a Latin square can be involved in at most n intercalates. Kotzig and Turgeon [9] proved that for all orders except 2×2 and 4×4 there exist Latin squares with no intercalates, and Heinrich and Wallis [4] proved that for all n there exist $n \times n$ Latin squares with $\Omega(n^3)$ intercalates.

In [12] McKay and Wanless conjectured the following.

*Department of Mathematics, ETH, 8092 Zürich. Email: matthew.kwan@math.ethz.ch.

[†]Department of Mathematics, ETH, 8092 Zürich. Email: benjamin.sudakov@math.ethz.ch. Research supported in part by SNSF grant 200021-149111.

Conjecture 1. *Let $\mathbf{N} = N(\mathbf{L})$ be the number of intercalates in a uniformly random Latin square $\mathbf{L} \in \mathcal{L}$. For any $\varepsilon > 0$, a.a.s.¹*

$$(1 - \varepsilon) \frac{n^2}{4} \leq \mathbf{N} \leq (1 + \varepsilon) \frac{n^2}{4}.$$

They were able to prove the substantially weaker lower bound that a.a.s. $\mathbf{N} \geq n^{3/2-\varepsilon}$ for any $\varepsilon > 0$. Before our paper, the best upper bound was due to Cavenagh, Greenhill and Wanless [2], who proved that a.a.s. $\mathbf{N} \leq (9/2)n^{5/2}$. The techniques used for these upper and lower bounds are very different, and we incorporate both to prove the following improved bounds. In particular we are able to prove the lower bound in Conjecture 1.

Theorem 1. *Let $\mathbf{N} = N(\mathbf{L})$ be the number of intercalates in a uniformly random Latin square $\mathbf{L} \in \mathcal{L}$. First,*

$$(1 - o(1)) \frac{n^2}{4} \leq \mathbb{E}\mathbf{N} \leq (1 + o(1)) \frac{n^2}{2}.$$

Second, for any fixed $\varepsilon > 0$ and any function $f \rightarrow \infty$, a.a.s.

$$(1 - \varepsilon) \frac{n^2}{4} \leq \mathbf{N} \leq fn^2.$$

Theorem 1 is an immediate corollary of two theorems that may be of independent interest, which we discuss in Section 2.

A different property that we believe typically holds for random Latin squares is that they have “low discrepancy” or are “quasirandom” in a certain sense. This is related to a conjecture by Linial and Luria [11]. To state their conjecture, note that \mathcal{L} can be interpreted as the set of all $n \times n \times n$ zero-one arrays with a single “1” in each axis-aligned line. To be specific, an $n \times n$ Latin square L corresponds to the $n \times n \times n$ array $A = A(L)$ where $A_{i,x,q} = 1$ if $L_{i,x} = q$. A *box* is a set of the form $T = I \times X \times Q$, where $I, X, Q \subseteq [n]$. For a box T , define its *volume* $\text{vol } T = |X||Y||Z|$. Let $N_T(L)$ be the number of ones in $A(L)$ in the positions in the box T . Linial and Luria’s conjecture is as follows.

Conjecture 2. *There exist arbitrarily large Latin squares L with the following property. For any box $T = I \times X \times Q$,*

$$\left| N_T(L) - \frac{\text{vol } T}{n} \right| = O(\sqrt{\text{vol } T}).$$

that is, they conjecture that there are Latin squares (zero-one arrays) such that in any box, the density of ones is very close to the density $1/n$ of ones in the entire $n \times n \times n$ array.

It is natural to expect that in fact the statement of Conjecture 2 holds a.a.s. for a uniformly random Latin square $\mathbf{L} \in \mathcal{L}$. Linial and Luria proved the weaker result that a.a.s. every empty box T with $N_T(\mathbf{L}) = 0$ has $\text{vol } T \leq n^2 \log^2 n$. We are able to show that random Latin squares a.a.s. have low discrepancy, although we are not able to reach the optimal conjectured discrepancy in Conjecture 2.

¹By “asymptotically almost surely”, or “a.a.s.”, we mean that the probability of an event is $1 - o(1)$. Here and for the rest of the paper, asymptotics are as $n \rightarrow \infty$

Theorem 2. *For a uniformly random $\mathbf{L} \in \mathcal{L}$, we a.a.s. have the following. For any box $T = I \times X \times Q$,*

$$\left| N_T(\mathbf{L}) - \frac{\text{vol } T}{n} \right| = O\left(\sqrt{\text{vol } T} \log n + n \log^2 n\right).$$

The proof of Theorem 2 is in Section 6.

2 Outline of the proof of Theorem 1

The first new ingredient for the proof of Theorem 1 is the following upper bound, both in expectation and with high probability, for the number of intercalates in two rows of a random Latin square.

Theorem 3. *Let $\mathbf{N}_2 = N_2(\mathbf{L})$ be the number of intercalates in the first two rows of a uniformly random Latin square $\mathbf{L} \in \mathcal{L}$. We have*

$$\mathbb{E}\mathbf{N}_2 \leq 1 + o(1)$$

and

$$\Pr(\mathbf{N}_2 \geq t) = e^{-\Omega(t \log t)}.$$

Note that $\mathbb{E}\mathbf{N} = \binom{n}{2} \mathbb{E}\mathbf{N}_2$ by linearity of expectation, so the upper bound on $\mathbb{E}\mathbf{N}$ in Theorem 1 immediately follows from Theorem 3. (Then, the a.a.s. upper bound on \mathbf{N} follows from Markov's inequality). We doubt that the bound on $\mathbb{E}\mathbf{N}_2$ in Theorem 3 is sharp; Conjecture 1 suggests that $\mathbb{E}\mathbf{N}_2 \sim 1/2$, and we expect that moreover \mathbf{N}_2 has an asymptotic Poisson distribution with this mean.

The second ingredient for Theorem 1 is a bound for the lower tail probability of the number of intercalates in a random Latin square.

Theorem 4. *There is a constant C such that the following holds. Let $\mathbf{N} = N(\mathbf{L})$ be the number of intercalates in a uniformly random Latin square $\mathbf{L} \in \mathcal{L}$. Suppose $\varepsilon \geq C \log^{1/3} n / n^{1/6}$. Then*

$$\Pr\left(\mathbb{E}\mathbf{N} < \frac{n^2}{4}(1 - \varepsilon)\right) \leq \exp(-\Omega(\varepsilon^2 \sqrt{n})).$$

Clearly Theorem 4 implies the a.a.s. lower bound on \mathbf{N} in Theorem 1, and because $\mathbf{N} \geq 0$ this in turn implies the lower bound on $\mathbb{E}\mathbf{N}$. We expect that Theorem 4 is far from optimal, and that the correct order of magnitude of the lower tail probability is $\exp(-\Omega(\varepsilon^2 n^2))$.

When studying random combinatorial structures with little independence, an indispensable technique is the analysis of “switching” operations that make local changes to an object. One defines switchings that affect some parameter in a controllable way, then estimates the number of ways to switch to and from each object, to understand the relative likelihood of each possible value of the parameter. Switchings underpin the proofs of both Theorem 3 and Theorem 4.

Latin squares are quite “rigid” objects, so one cannot easily define switching operations that make only small changes to a Latin square. In [2], Cavenagh, Greenhill and Wanless managed to overcome

this difficulty when studying two fixed rows of a random Latin square. They considered switchings that make wide-ranging, complicated changes to the whole Latin square, but have a controllable effect on the two rows of interest. We prove Theorem 3 with a simpler switching operation in a similar spirit. The details are in Section 3.

To prove Theorem 4, we use Theorem 3 and some ideas from [12]. A $k \times n$ *Latin rectangle* is a $k \times n$ array of the numbers from 1 to n , where each number appears once in each row and not more than once in each column. We denote the set of all $k \times n$ Latin rectangles by \mathcal{L}_k .

For $k \leq n$, any $k \times n$ Latin rectangle can be extended to a $n \times n$ Latin square. The number of ways to do this does not depend too much on the Latin rectangle. Indeed, for a $k \times n$ Latin rectangle $L \in \mathcal{L}_k$ let $\mathcal{L}^*(L) \subseteq \mathcal{L}$ be the set of $n \times n$ Latin squares whose first k rows coincide with L . The following estimate is proved with standard upper and lower bounds on the permanent, and is proved in Section 7.

Proposition 5. *For Latin rectangles $L, L' \in \mathcal{L}_k$,*

$$\frac{|\mathcal{L}^*(L)|}{|\mathcal{L}^*(L')|} \leq e^{O(n \log^2 n)},$$

uniformly over k .

So, the strategy is to find a lower bound on the number of intercalates in a random $k \times n$ Latin rectangle (for some k to be determined) that holds with very high probability. We will then be able to apply Proposition 5 to show that the number of intercalates in the first k rows of a random Latin square satisfies the same bound with high probability. We can use the union bound to show that this holds simultaneously for many choices of k rows, which gives a lower bound for the total number of intercalates in a random Latin square.

In [12], McKay and Wanless studied the number of intercalates in a random Latin rectangle. Using their methods, we will prove the following estimate.

Lemma 6. *There is a constant C such that the following holds. Let $\mathbf{L} \in \mathcal{L}_k$ be a uniformly random $k \times n$ Latin rectangle, conditioned on the event that no row is involved in more than K intercalates. Suppose $k \geq \sqrt{n}$ and $k \leq K \leq n$.*

Let $\mathbf{N} = N(\mathbf{L})$ be the number of intercalates in \mathbf{L} . If $t \geq Ck^2(K/n + k/K)$, then

$$\Pr\left(\mathbf{N} \leq \frac{k^2}{4} - t\right) \leq e^{-\Omega(t^2/k^2)}.$$

Note that Theorem 3 and the union bound imply that with high probability no row is involved in many intercalates, which will give us an appropriate value of K with which to apply Lemma 6. In Section 4 we prove Lemma 6, and in Section 5 we give the details of how to combine Proposition 5, Lemma 6 and Theorem 3 to obtain Theorem 4.

3 Proof of Theorem 3

Note that any two rows i, j of a Latin square $L \in \mathcal{L}$ define a permutation $\sigma_{i,j}(L)$ on the columns of L : column x maps to the column y with $L_{i,y} = L_{j,x}$. Note that this permutation is a derangement;

$L =$	1	3	5	4	2
	4	5	3	2	1
	2		4		
	3		2		
	5		1		

$\text{fx}_{\{2,3\}}(L) =$	1	5	3	4	2
	4	3	5	2	1
	2		4		
	3		2		
	5		1		

Figure 1. We show the effect of the fix operation. (Not all cells are depicted). Note that $\{1, 3\}$ was not flippable in L but is flippable in $\text{fx}_{\{2,3\}}(L)$. Also note that $\sigma_{1,2}(L) = \sigma_{1,2}(\text{fx}_{\{2,3\}}(L)) = (145)(23)$.

$L =$	1	3	5	4	2
	4	5	3	2	1
	2		4		
	3		2		
	5		1		

$\text{fl}_{\{1,3\}}(L) =$	1	3	5	4	2
	3	5	4	2	1
	4		2		
	2		3		
	5		1		

Figure 2. We show the effect of the flip operation. Note that $\sigma_{1,2}(L) = (145)(23)$ and $\sigma_{1,2}(\text{fl}_{\{1,3\}}(L)) = (12345)$.

it has no fixed points $x \mapsto x$. We will be concerned with the permutation $\sigma_{1,2}(L)$ defined by the first two rows. This permutation decomposes into cycles, the set of which we denote $C(L)$. (For our purposes a cycle is a set of columns). Let $C^\alpha(L) \subseteq C(L)$ be the set of cycles of length α , and for a column x let $c_x(L) \in C(L)$ be the cycle which contains x .

Let $N^\alpha(L) = |C^\alpha(L)|$. We are interested in $\mathbf{N}_2 = N_2(L) = N^2(L)$ for a uniformly random $L \in \mathcal{L}$. Let $\mathcal{L}(s) \subseteq \mathcal{L}$ be the set of Latin squares L with $N^2(L) = s$.

We first define two primitive switching operations on a Latin square, chosen such that they have a controllable effect on the intercalate count in the first two rows.

Definition 7. Consider a Latin square L .

- For any cycle $c \in C(L)$, we can obtain a new Latin square $\text{fx}_c(L)$ by exchanging the contents of rows 1 and 2, for each column in c . We call this operation *fix*. We also write $\text{fx}_x(L)$ to denote $\text{fx}_{c_x(L)}(L)$.
- Just as two rows i, j of L define the permutation $\sigma_{i,j}(L)$, every two columns x and y also define a permutation $\tau_{x,y}(L)$ of rows (with row i mapping to the row j which satisfies $L_{j,x} = L_{i,y}$). In the cycle decomposition of $\tau_{x,y}(L)$, if rows 1 and 2 are in different cycles c_1 and c_2 , then we say $\{x, y\}$ is a *flippable* pair, and write $\{x, y\} \in \text{FL}(L)$. We can obtain a new Latin square $\text{fl}_{\{x,y\}}(L)$ by exchanging column x and y for each row in c_2 . We call this operation *flip*.

We will be using the flip operation to merge two cycles into a larger cycle, and we will be using the fix operation to make a pair of columns from different cycles flippable, if necessary. To justify this, we make a number of simple observations about the properties of the fix and flip operations.

Fact 8. The operations in Definition 7 have the following consequences.

1. Suppose $\{x, y\} \in \text{FL}(L)$. If $c_x(L) \neq c_y(L)$, with say $c_x(L) \in C^\alpha(L)$ and $c_y(L) \in C^\beta(L)$. Let $L' = \text{fl}_{\{x,y\}}(L)$. Then

$$c_x(L') = c_y(L') = c_x(L) \cup c_y(L) \in C^{\alpha+\beta}(L').$$

Also, $C(L) \setminus \{c_x(L), c_y(L)\} = C(L') \setminus \{c_x(L')\}$. That is, flipping with x and y merges $c_x(L)$ and $c_y(L)$ and leaves the other cycles unaffected.

2. If $c \in C^2(L)$ is an intercalate, then $\sigma_{1,2}(L) = \sigma_{1,2}(\text{fx}_c(L))$. That is, the fix operation does not change the induced permutation.
3. Suppose $c_x(L) \neq c_y(L)$ and $\{x, y\} \in \text{FL}(L)$ (respectively, $\{x, y\} \notin \text{FL}(L)$). Let $L' = \text{fx}_x(L)$. Then $\{x, y\} \notin \text{FL}(L')$ (respectively $\{x, y\} \in \text{FL}(L')$). That is, the fix operation changes the flippability of $\{x, y\}$.
4. For any cycle c , we have $\text{fx}_c(i)(\text{fx}_c(L)) = L$. For any $\{x, y\} \in \text{FL}(L)$, we have $\text{fl}_{\{x, y\}}(\text{fl}_{\{x, y\}}(L)) = L$. That is, the fix and flip operations are both involutions and in particular they are invertible.
5. Suppose $\{x, y\} \in \text{FL}(L)$ and $c_x(L) \neq c_y(L)$, with $c_y(L) \in C^2(L)$. Let $\sigma' = \sigma_{1,2}(\text{fl}_{\{x, y\}}(L))$. Then $(\sigma')^2(x) = y$.

With these observations in mind we can define a compound operation that merges an intercalate with another cycle, regardless of flippability.

Definition 9. For columns x, y with $c_x(L) \neq c_y(L)$ and $c_y(L) \in C^2(L)$ define

$$\text{jo}_{x,y}(L) = \begin{cases} \text{fl}_{\{x,y\}}(L) & \text{if } \{x, y\} \in \text{FL}(L), \\ \text{fl}_{\{x,y\}}(\text{fx}_y(L)) & \text{if } \{x, y\} \notin \text{FL}(L). \end{cases}$$

We call this operation *join*. If also $c_x(L) \in C^2(L)$ this is a *double* join, otherwise it is a *single* join.

Let $\mathcal{L}(s) \subseteq \mathcal{L}$ be the set of Latin squares L with s intercalates in the first two rows (that is, with $N^2(L) = s$).

We make some observations about the join operation.

Fact 10. Single and double joins have the following consequences.

1. A single join always decreases $N^2(\cdot)$ by exactly one, and the merged cycle has length greater than 4.
2. A double join always decreases $N^2(\cdot)$ by exactly two, and the merged cycle has length 4.
3. For $L \in \mathcal{L}(s+1)$, the number of Latin squares $L' \in \mathcal{L}(s)$ which we can reach with a single join is

$$(n - 2N^2(L)) \times 2N^2(L) = 2(s+1)(n - 2(s+1)).$$

(Choose a column x not in an intercalate and a column y in an intercalate).

4. For $L \in \mathcal{L}(s+2)$, the number of Latin squares $L' \in \mathcal{L}(s)$ which we can reach with a double join is

$$2N^2(L) \times 2(N^2(L) - 1) = 4(s^2 + 3s + 2) \geq 2s^2.$$

(Choose a column x in an intercalate and a column y in a different intercalate).

5. For $L' \in \mathcal{L}(s)$, the number of Latin squares $L \in \mathcal{L}(s+1)$ which can reach L' with a single join is at most

$$2(n - 2N^2(L) - 3N^3(L) - 4N^4(L)) \leq 2(n - 2s).$$

(For $\sigma' = \sigma_{1,2}(L')$, choose a column x in a cycle with length greater than 4, and let $y = (\sigma')^2(x)$. If $\{x, y\} \in \text{FL}(L')$ then flip, and then either fix or don't).

6. For $L' \in \mathcal{L}(s)$, the number of Latin squares $L \in \mathcal{L}(s+2)$ which can reach L' with a double join is at most

$$2 \times 4N^4(L) \leq 2n.$$

(For $\sigma' = \sigma_{1,2}(L')$, choose a column x in a 4-cycle, let $y = (\sigma')^2(x)$, flip if possible and then either fix or don't).

Let $J(s)$ be the number of ways to single join from a Latin square in $\mathcal{L}(s+1)$ to one in $\mathcal{L}(s)$. We have

$$\begin{aligned} 2(s+1)(n-2(s+1))|\mathcal{L}(s+1)| &= J(s) \leq 2(n-2s)|\mathcal{L}(s)|, \\ \frac{|\mathcal{L}(s+1)|}{|\mathcal{L}(s)|} &\leq \frac{n-2s}{(s+1)(n-2s-2)}. \end{aligned}$$

Similarly, double-counting the number of ways to double join from a Latin square in $\mathcal{L}(s+2)$ to one in $\mathcal{L}(s)$, we obtain

$$\frac{|\mathcal{L}(s+2)|}{|\mathcal{L}(s)|} \leq \frac{n}{s^2}.$$

So,

$$\frac{|\mathcal{L}(s+1)|}{|\mathcal{L}(s)|} \leq \frac{1}{s+1} \left(1 + O\left(\frac{1}{n}\right) \right)$$

for $2s \leq n/2$ and $|\mathcal{L}(s+2)|/|\mathcal{L}(s)| \leq 1$ for $2s \geq n/2$ (for large n). It follows that

$$\Pr(\mathbf{N}_2 = t) \leq \frac{|\mathcal{L}(t)|}{|\mathcal{L}(0)|} \leq \prod_{s=0}^{t-1} \frac{|\mathcal{L}(s+1)|}{|\mathcal{L}(s)|} = \frac{1}{t!} e^{O(t/n)}$$

and

$$\Pr\left(t < \mathbf{N}_2 \leq \frac{n}{4}\right) \leq O(1) \sum_{s=t}^{n/2} \frac{1}{s!} \leq O\left(\frac{1}{t!}\right) \sum_{r=0}^{\infty} \frac{1}{t^r} = O\left(\frac{1}{t!}\right) = e^{-\Omega(t \log t)}$$

for $t \leq n/4$, and

$$\begin{aligned} \Pr(\mathbf{N}_2 = t) &= O\left(\frac{1}{(n/4)!}\right), \\ \Pr(\mathbf{N}_2 > t) &\leq O\left(\frac{n}{(n/4)!}\right) = e^{-\Omega(n \log n)} = e^{-\Omega(t \log t)} \end{aligned}$$

for $t > n/4$. It follows that $\Pr(\mathbf{N}_2 \geq t) = e^{-\Omega(t \log t)}$ for all t .

We now bound $\mathbb{E}\mathbf{N}_2$. Fix some constant C , fix some k with $1 \ll k \ll n$, and define each d_t ($t \leq k$) by

$$\frac{\Pr(\mathbf{N}_2 = t)}{\Pr(\mathbf{N}_2 = 0)} = \frac{|\mathcal{L}(t)|}{|\mathcal{L}(0)|} = \frac{1}{t!} e^{Ck/n} - d_t.$$

$$L = \begin{array}{|c|c|c|c|c|c|} \hline 1 & 3 & 5 & 4 & 2 & 6 \\ \hline & & 2 & 3 & & \\ \hline \end{array} \quad \text{tw}_{\{(3,1,2),(4,6,5)\}}^1(L) = \begin{array}{|c|c|c|c|c|c|} \hline 5 & 1 & 3 & 2 & 6 & 4 \\ \hline & & 2 & 3 & & \\ \hline \end{array}$$

Figure 3. We show the effect of the twist operation to create an intercalate involving (1, 3) and (1, 4).

For large C each d_t is positive. We have

$$1 \sim \sum_{t=0}^k \Pr(\mathbf{N}_2 = t) \sim \Pr(\mathbf{N}_2 = 0) \left(e^{Ck/n} \sum_{t=0}^k 1/t! - \sum_{t=0}^k d_t \right) \sim \Pr(\mathbf{N}_2 = 0) \left(e - \sum_{t=1}^k d_t \right)$$

(note that $d_0 = o(1)$) and

$$\mathbb{E} \mathbf{N}_2 = \sum_{t=0}^k t \Pr(\mathbf{N}_2 = t) = \Pr(\mathbf{N}_2 = 0) \left(e^{Ck/n} \sum_{t=0}^k 1/(t-1)! - \sum_{t=0}^k t d_t \right) \sim \frac{e - \sum_{t=1}^k t d_t}{e - \sum_{t=1}^k d_t} \leq 1.$$

4 Proof of Lemma 6

Let $\mathcal{L}_k^K \subseteq \mathcal{L}_k$ be the set of Latin rectangles L in which no row is involved in more than K intercalates. (We say these Latin rectangles are “good”). Let $\mathcal{L}_k^K(s) \subseteq \mathcal{L}_k^K$ be the set of good Latin rectangles with exactly s intercalates.

Definition 11. Consider a Latin rectangle $L \in \mathcal{L}_k^K$. For a row i and a cyclically ordered set of columns $(x y z)$, we obtain a new $k \times n$ array $L' = \text{ro}_{(xyz)}^i(L)$ by swapping the symbols in positions (i, x) , (i, y) , (i, z) in a cyclic fashion: $L'_{i,x} = L_{i,z}$, $L'_{i,y} = L_{i,x}$, $L'_{i,z} = L_{i,y}$. This operation is called *rotate*. Note that L' might not be a Latin rectangle, because we might have caused a column to contain two of the same symbol.

Now, we define the *twist* operation. For a Latin rectangle $L \in \mathcal{L}_k^K$, a row i and distinct columns x, y, z, x', y', z' , let $L' = \text{ro}_{(x' y' z')}^i(\text{ro}_{(xyz)}^i(L))$. Suppose the following conditions are satisfied.

- The rectangle L' is a Latin rectangle, and it is good (that is, $L' \in \mathcal{L}_k^K$).
- The positions (i, y) , (i, z) , (i, y') , (i, z') are involved in no intercalates in L or in L'
- The positions (i, x) and (i, x') are involved in no intercalates in L , and in L' there is an intercalate involving both (i, x) and (i, x') .

Then we define the twist of L by $\text{tw}_{\{(x,y,z),(x',y',z')\}}^i(L) = L'$. Note that we might have created (up to $O(k)$) intercalates involving (i, x) and (i, x') , other than the “main” intercalate involving both (i, x) and (i, x') . If the only new intercalate in L' is the main one, then we say the twist was *simple*.

Lemma 12. *The number of good Latin rectangles $L' \in \mathcal{L}_k^K(s+1)$ which we can reach via a simple twist from a specific good Latin rectangle $L \in \mathcal{L}_k^K(s)$ is at least*

$$\frac{1}{2} k^2 n^4 \left(1 - O\left(\frac{1}{k} + \frac{k}{n} + \frac{K}{n} + \frac{s}{kK} \right) \right).$$

Proof. Let $\Psi(L)$ be the set of rows of L involved in exactly K intercalates. We have $|\Psi(L)| \leq 4s/K$. Now, choose rows i and j not in $\Psi(L)$, in which we will create an intercalate. There are

$$\left(k - \frac{4s}{K}\right) \left(k - \frac{4s}{K} - 1\right) = k^2 \left(1 - O\left(\frac{s}{kK} + \frac{1}{k}\right)\right)$$

ways to do this. Next, choose distinct columns x, y, x', y' . To create an intercalate in columns x and x' , let z' be the unique column with $L_{j,x} = L_{i,z'}$, and let z be the column with $L_{j,x'} = L_{i,z}$. There are $n^4(1 + O(1/n))$ ways to make these choices, but some of these do not give rise to a valid twist operation. Let $L' = \text{ro}_{(xyz)}^i(\text{ro}_{(x'y'z')}^i(L))$; the possible violations are as follows.

- The symbol $L_{i,x}$ might already appear in column y (so that L' is not a Latin rectangle). For any x', y, y' there are at most k choices of x with this property, so we should subtract kn^3 for our upper bound. Similarly $L_{i,x'}$, $L_{i,y}$, $L_{i,y'}$, $L_{i,z}$ or $L_{i,z'}$ might appear in column y' , z , z' , x or x' respectively. We should therefore subtract $6kn^3$.
- We might have $z' \in \{x', y, y', z\}$. For any x', y, y' there are at most 4 choices of x that cause this. Similarly we might have $z \in \{x, y, y', z'\}$. We should subtract $8n^3$ to compensate for both.
- One of the positions (i, x) , (i, x') , (i, y) , (i, y') , (i, z) or (i, z') might already be involved in an intercalate. We should subtract $6Kn^3$ to compensate for this.
- There might be an intercalate involving (i, y) in L' . This can only occur if for one of the $k-1$ columns (w , say) in L' which contains the symbol $L_{i,x}$ (in row $q \neq 1$, say), we have $L'_{i,w} = L'_{q,y}$. For any x, x', y' there are at most k choices of y for which this occurs. Similarly, putting $L_{i,y}$, $L_{i,z}$, $L_{i,x'}$, $L_{i,y'}$ or $L_{i,z'}$ in position (i, z) , (i, x) , (i, y') , (i, z') or (i, x') respectively might create an intercalate involving that position (other than the one given by positions (i, x) , (i, x') , (j, x) , (j, x')). So we should subtract $6kn^3$ to compensate for this.

If the above points are not violated then we can use i, x, y, z, x', y', z' to do a simple twist, so the number of valid ways to simple twist is at least

$$\begin{aligned} & \frac{1}{2}k^2 \left(1 - O\left(\frac{s}{kK} + \frac{1}{k}\right)\right) \left(n^4 \left(1 + O\left(\frac{1}{n}\right)\right) - O(Kn^3) - O(kn^3) - O(n^3)\right) \\ &= \frac{1}{2}k^2 n^4 \left(1 - O\left(\frac{1}{k} + \frac{k}{n} + \frac{K}{n} + \frac{s}{kK}\right)\right). \end{aligned}$$

(We divide by 2 to compensate for the fact that we can exchange x and x' , y and y' , z and z' to give the same twist). \square

Lemma 13. *The number of good Latin rectangles $L \in \mathcal{L}_k^K$ from which we can twist to a specific good Latin rectangle $L' \in \mathcal{L}_k^K(s)$ is at most $2sn^4$.*

Proof. Twisting from L must have created one of the s intercalates in L' as its main intercalate, operating in one of its two rows. The columns $\{x, x'\}$ are determined by the intercalate that was created, and there are at most n^4 choices of y, y', z, z' that could have been used. So the number of Latin rectangles L that can twist to L' is at most $2sn^4$. \square

We can use Lemmas 12 and 13 to give an upper and lower bound on the number of ways to simple twist from a Latin rectangle in $\mathcal{L}_k^K(s-1)$ to a Latin rectangle in $\mathcal{L}_k^K(s)$. For $s = O(k^2)$, $k \geq \sqrt{n}$ and $k \leq K \leq n$, this gives

$$\frac{|\mathcal{L}_k^K(s-1)|}{|\mathcal{L}_k^K(s)|} \leq \frac{s}{k^2/4} \exp\left(O\left(\frac{K}{n} + \frac{k}{K}\right)\right),$$

so for $t \geq 0$,

$$\frac{|\mathcal{L}_k^K(k^2/4 - s)|}{|\mathcal{L}_k^K(k^2/4)|} \leq \prod_{r=0}^{s-1} \frac{|\mathcal{L}_k^K(k^2/4 - r - 1)|}{|\mathcal{L}_k^K(k^2/4 - r)|} \leq \prod_{r=0}^{s-1} \left(\frac{k^2/4 - r}{k^2/4} \exp\left(O\left(\frac{K}{n} + \frac{k}{K}\right)\right) \right).$$

For $0 \leq 2r \leq k^2/4$ we have $(k^2/4 - r)/(k^2/4) = \exp(-\Theta(r/k^2))$; it follows that

$$\begin{aligned} \Pr(\mathbf{N} = k^2/4 - s) &\leq \frac{|\mathcal{L}_k^K(k^2/4 - s)|}{|\mathcal{L}_k^K(k^2/4)|} \\ &\leq \exp\left(-\left(\sum_{r=0}^{s-1} \Theta\left(\frac{r}{k^2}\right)\right) + O\left(s\left(\frac{K}{n} + \frac{k}{K}\right)\right)\right) \\ &= \exp\left(-\Omega\left(\frac{s^2}{k^2}\right) + O\left(s\left(\frac{K}{n} + \frac{k}{K}\right)\right)\right). \end{aligned}$$

If $t \geq Ck^2(K/n + k/K)$ for large C , then

$$\Pr(\mathbf{N} < k^2/4 - t) \leq \sum_{s=t}^{k^2/4} e^{-\Omega(s^2/k^2)} = e^{-\Omega(t^2/k^2)}.$$

5 Synthesis: Proof of Theorem 4

In this proof, the notation $a \ll b$ means that b/a is sufficiently large (contrary to the usual meaning that $b/a \rightarrow \infty$).

Suppose k and K satisfy $\sqrt{n} \log n / \varepsilon \ll k \ll n / \log n$, $k \log n \ll K \leq n$ and $\varepsilon \gg K/n + k/K$. Let \mathcal{E} be the event that none of the first k rows of \mathbf{L} are involved in more than K intercalates, in the Latin rectangle induced by the first k rows. By Theorem 3 and the union bound (and symmetry considerations),

$$\Pr(\mathcal{E}^c) = \exp\left(-\Omega\left(\frac{K}{k} \log \frac{K}{k}\right)\right).$$

Let \mathbf{N}_k be the number of intercalates in the first k rows of \mathbf{L} . For any fixed $\varepsilon > 0$, by Lemma 6 and Proposition 5,

$$\Pr(\mathbf{N}_k < (1 - \varepsilon/2)k^2/4 \mid \mathcal{E}) = e^{-\Omega(\varepsilon^2 k^2)} e^{O(n \log^2 n)} = e^{-\Omega(\varepsilon^2 k^2)}$$

We are assuming that $\varepsilon k \gg \sqrt{n} \log n$ so $\varepsilon^2 k^2 \gg (K/k) \log(K/k)$ and it follows that

$$\Pr(\mathbf{N}_k < (1 - \varepsilon/2)k^2/4) \leq \exp\left(-\Omega\left(\frac{K}{k} \log \frac{K}{k}\right)\right) \tag{1}$$

unconditionally. Now, fix some $M \in \mathbb{N}$ that is no more than say half of the reciprocal of the probability in (1). Choose M independent uniformly random sets of k rows. For a given pair of rows, the probability a uniformly random set of k rows contains that pair is $p = \binom{n}{k-2} / \binom{n}{k} = (k/n)^2 (1 + O(\frac{k}{n}))$. By the Chernoff bound and the union bound, a.a.s. every pair is contained in

$$Mp + O(\sqrt{Mp} \log n) = Mp \exp\left(1 + \frac{\log n}{\sqrt{Mp}}\right)$$

of our M sets of k rows. Fix such a collection of sets of rows. By the union bound and symmetry, except with the probability in (1) each of our sets of k rows contains at least $(1 - \varepsilon)k^2/4$ intercalates, so assuming $k/n, \log n/\sqrt{Mp} \ll \varepsilon$, we have

$$\begin{aligned} N &\geq \frac{M(1 - \varepsilon/2)k^2/4}{Mp(1 + \log n/\sqrt{Mp})} \\ &\geq (1 - \varepsilon) \frac{n^2}{4}. \end{aligned}$$

It remains to optimize for k, K and M (and determine for which ε feasible values for k, K, M exist). We find that we can choose $k \asymp \sqrt{n} \log n / \varepsilon$ and $K \asymp \varepsilon n$, yielding a tail bound of $\exp(-\Omega(\varepsilon^2 \sqrt{n}))$, provided $\varepsilon \gg \log^{1/3} n / n^{1/6}$.

6 Proof of Theorem 2

Fix a box $T = I \times X \times Q$ (there are $(2^n)^3 = 8^n$ possible choices). We will show that the bound on $N_T(\mathbf{L})$ in Theorem 2 holds with probability $o(8^{-n})$, which will allow us to apply the union bound over choices of T .

For a Latin square L , we define a bipartite graph $G_Q(L)$ as follows. Both parts have n vertices (we abuse notation and say the vertex set is $[n] \sqcup [n]$); one of the parts is identified with the set of rows of the Latin square and the other part is identified with the set of columns. For each row i and column x such that $L_{i,x} \in Q$, we put an edge between i and x in $G_Q(L)$. Now, the number of ones $N_T(L)$ in T is just the number of edges $e_{G_Q(L)}(I, X)$ between I and X in $G_Q(L)$.

Let \mathcal{G}_d be the set of d -regular bipartite graphs on $[n] \sqcup [n]$. For $G \in \mathcal{G}_{|Q|}$, let $|\mathcal{L}^*(G)|$ be the number of Latin squares L with $G_Q(L) = G$. We can use standard bounds on the permanent to prove that $|\mathcal{L}^*(G)|$ does not vary very much with G . See Section 7 for details.

Proposition 14. *For a set of symbols Q and $|Q|$ -regular bipartite graphs G and G' ,*

$$\frac{|\mathcal{L}^*(G)|}{|\mathcal{L}^*(G')|} \leq e^{O(n \log^2 n)}$$

uniformly over Q .

So, $G_Q(L)$ is not too far from the uniform distribution on $\mathcal{G}_{|Q|}$, and events that hold with very high probability for a uniformly random $G \in \mathcal{G}_{|Q|}$ also hold with very high probability for $G_Q(L)$.

It is possible to obtain discrepancy tail bounds for random regular (bipartite) graphs using switchings of the type in [10, Theorem 2.2]. Such a bound would nearly provide the result we are after (although

there would be difficulties for very dense graphs). However, at the range of probabilities we are interested in, regular bipartite graphs comprise a non-negligible proportion of all bipartite graphs with the appropriate number of edges, and (modulo a deep enumeration theorem for regular bipartite graphs) this enables a simpler approach. Let $\mathbb{B}(n, p)$ be the random graph distribution on the vertex set $[n] \sqcup [n]$, where each of the n^2 possible edges between the parts are present with independent probability p .

Lemma 15. *Let $p = d/n$. The probability a random graph $\mathbf{B} \in \mathbb{B}(n, p)$ is d -regular is $e^{-O(n \log n)}$. Also, conditioning on this event gives the uniform distribution on \mathcal{G}_d .*

To prove Lemma 15 we will use the following theorem of Wormald and Liebenau [?].

Theorem 16. *Let $p = d/n$. The number of d -regular bicoloured graphs with $2n$ vertices is asymptotic to*

$$\left(p^p(1-p)^{1-p}\right)^{n^2} \binom{n}{pn}^{2n} n \sqrt{2\pi p(1-p)} e^{-1/2}.$$

Proof of Lemma 15. The probability \mathbf{B} has exactly $dn = pn^2$ edges is

$$\binom{n^2}{pn^2} p^{pn^2} (1-p)^{pn^2} \asymp \frac{1}{n \sqrt{p(1-p)}}.$$

(here we used Stirling's approximation). By symmetry, each graph with dn edges is equally likely. By Theorem 16, the fraction of such graphs which are d -regular is

$$\begin{aligned} & \left(p^p(1-p)^{1-p}\right)^{n^2} \binom{n}{pn}^{2n} n \sqrt{2\pi p(1-p)} e^{-1/2} \bigg/ \binom{n^2}{pn^2} \\ &= (p(1-p)n)^{-(2+o(1))n} \leq e^{-O(n \log n)}. \end{aligned} \quad \square$$

Now, discrepancy in $\mathbb{B}(n, p)$ (for $p = |Q|/n$) is very easy to study. Indeed, for $\mathbf{B} \in \mathbb{B}(n, p)$ the law of $e_{\mathbf{B}}(I, X)$ is the binomial distribution $\text{Bin}(|I||X|, p)$ with mean $|I||X|p = \text{vol } T/n$. Let $\mathbf{G} \in \mathcal{G}_{|Q|}$ be a uniformly random $|Q|$ -regular bipartite graph. By a binomial large deviation inequality (for example [6, Theorem 2.1]), Proposition 14 and Lemma 15, we have

$$\begin{aligned} \Pr\left(\left|N_T(\mathbf{L}) - \frac{\text{vol } T}{n}\right| > t\right) &= \Pr\left(\left|e_{G_Q(\mathbf{L})}(I, X) - \frac{\text{vol } T}{n}\right| > t\right) \\ &\leq \Pr\left(\left|e_{\mathbf{G}}(I, X) - \frac{\text{vol } T}{n}\right| > t\right) e^{O(n \log^2 n)} \\ &\leq \Pr\left(\left|e_{\mathbf{B}}(I, X) - \frac{\text{vol } T}{n}\right| > t\right) e^{O(n \log^2 n + n \log n)} \\ &= \exp\left(-\Omega\left(\frac{t^2}{\text{vol } T/n + t}\right) + O(n \log^2 n)\right). \end{aligned}$$

If t is a large multiple of $\sqrt{\text{vol } T} \log n + n \log^2 n$, then this probability is $e^{-\Omega(n \log^2 n)} = o(8^{-n})$.

7 Proofs of Proposition 5 and Proposition 14

We can interpret a Latin square as a 1-factorization of $K_{n,n}$ (that is, an edge colouring with n colours). The correspondence is that the edge between vertex i in the first part and vertex x in the second part receives colour q if $L_{i,x} = q$.

Let $\Phi(G)$ be the number of 1-factorizations of a graph G . Schrijver [14] proved a lower bound on $\Phi(G)$ for regular bipartite graphs, and we can obtain an upper bound by iterating Brégman's theorem [1]. Both bounds are summarized in the following theorem.

Theorem 17. *Let G be a d -regular bipartite graph on $n + n$ vertices. Then*

$$\left(\frac{d!^2}{d^d}\right)^n \leq \Phi(G) \leq \prod_{k=1}^d (k!)^{n/k}.$$

Stirling's approximation gives $\Phi(G) \geq \left(\Omega\left(d(d/e^2)^d\right)\right)^n$, and additionally using the approximation $H_k := \sum_{i=1}^k \frac{1}{i} = \Theta(\log k)$ for the harmonic series,

$$\Phi(G) \leq \prod_{k=1}^d \left(O(\sqrt{k})\right)^{n/k} \left(\frac{k}{e}\right)^n \leq d^{O(n \log d)} \left(\frac{d!}{e^d}\right)^n \leq e^{O(n \log^2 d)} \left(O\left(\sqrt{d}(d/e^2)^d\right)\right)^n.$$

We have the following as a consequence.

Proposition 18. *Let G and G' be d -regular bipartite graphs on $n + n$ vertices. Then*

$$\frac{\Phi(G)}{\Phi(G')} \leq e^{O(n \log^2 d)}.$$

Now, in the notation of Proposition 14, note that $\mathcal{L}^*(G)$ is the number of 1-factorizations of G , times the number of 1-factorizations of its complement. For Proposition 5, we need to exchange the role of symbols and rows in our correspondence between Latin squares and 1-factorizations of $K_{n,n}$: put an edge coloured i between x and q if $L_{i,x} = q$. Then, $\mathcal{L}^*(L)$ is the number of 1-factorizations in the complement of a certain d -regular bipartite graph, which is an $(n - d)$ -regular bipartite graph. Therefore both Proposition 14 and Proposition 5 follow from Proposition 18.

8 Concluding remarks

We have shown that the number of intercalates \mathbf{N} in a uniformly random $n \times n$ Latin square a.a.s. satisfies $(1 - o(1))n^2/4 \leq \mathbf{N} \leq fn^2$, for any $f \rightarrow \infty$, and we showed that $(1 + o(1))n^2/4 \leq \mathbb{E}\mathbf{N} \leq (1 + o(1))n^2/2$. In doing so we obtained an exponentially-decaying estimate for the lower tail of \mathbf{N} and an exponential upper-tail estimate for the number of intercalates in two fixed rows. We also proved that random Latin squares typically have low discrepancy.

There are a number of related problems that remain open. First, there is the task of reducing the a.a.s. upper bound on \mathbf{N} to $(1 + o(1))n^2/4$ or at least to $O(n^2)$. The most obvious way of

approaching this would be to imitate our proof of the lower bound, and show that for some k satisfying $\sqrt{n} \log n \ll k$, with very high probability a random $k \times n$ Latin rectangle does not have too many intercalates. The tools from [12] can accomplish this conditioned on the nonexistence of “problematic configurations” of intercalates, but showing these configurations are unlikely appears to be a surprisingly difficult task.

Second, there is the problem of understanding the existence and number of substructures other than intercalates in random Latin squares. McKay and Wanless [12] conjecture that the number of 3×3 Latin subsquares should have expectation $\Theta(1)$, and similar logic would suggest that a.a.s. there are no Latin subsquares of larger order. A proof of either of these facts would be interesting. One can also ask about Latin subrectangles, or partial Latin subrectangles, of various kinds.

Third, there is the task of making further progress towards Conjecture 2. Even a slight improvement over our Theorem 2 would be interesting, because such an improvement would have to overcome the error introduced by theorems of the type in Section 7.

Finally, one might hope to prove results of the type in this paper for more general types of random designs, such as Latin cubes or Steiner triple systems. Unfortunately, results of the type in Section 7 are not readily available in these cases, which considerably limits the tools available. The breakthrough methods of Keevash [7, 8] for completing quasirandom partial designs may be useful.

References

- [1] LM Brgman. Certain properties of nonnegative matrices and their permanents. In *Dokl. Akad. Nauk SSSR*, volume 211, page 2730, 1973.
- [2] Nicholas J Cavenagh, Catherine Greenhill, and Ian M Wanless. The cycle structure of two rows in a random Latin square. *Random Structures & Algorithms*, 33(3):286–309, 2008.
- [3] Charles J Colbourn and Jeffrey H Dinitz. *Handbook of combinatorial designs*. CRC press, 2006.
- [4] Katherine Heinrich and WD Wallis. The maximum number of intercalates in a latin square. In *Combinatorial Mathematics VIII*, pages 221–233. Springer, 1981.
- [5] Mark T Jacobson and Peter Matthews. Generating uniformly distributed random Latin squares. *Journal of Combinatorial Designs*, 4(6):405–437, 1996.
- [6] Svante Janson, Tomasz Łuczak, and Andrzej Ruciński. *Random Graphs*. Cambridge University Press, 2000.
- [7] Peter Keevash. The existence of designs. *arXiv preprint arXiv:1401.3665*, 2014.
- [8] Peter Keevash. Counting designs. *arXiv preprint arXiv:1504.02909*, 2015.
- [9] Anton Kotzig and Jean Turgeon. On certain constructions for latin squares with no latin subsquares of order two. *Discrete Mathematics*, 16(3):263–270, 1976.
- [10] Michael Krivelevich, Benny Sudakov, Van H Vu, and Nicholas C Wormald. Random regular graphs of high degree. *Random Structures & Algorithms*, 18(4):346–363, 2001.

- [11] Nathan Linial and Zur Luria. Discrepancy of high-dimensional permutations. *arXiv preprint arXiv:1512.04123*, 2015.
- [12] Brendan D McKay and Ian M Wanless. Most Latin squares have many subsquares. *Journal of Combinatorial Theory, Series A*, 86(2):323–347, 1999.
- [13] Arthur O Pittenger. Mappings of Latin squares. *Linear algebra and its applications*, 261(1):251–268, 1997.
- [14] Alexander Schrijver. Counting 1-factors in regular bipartite graphs. *Journal of Combinatorial Theory, Series B*, 72(1):122–135, 1998.