RESTRICTED SET ADDITION IN GROUPS, II. A GENERALIZATION OF THE ERDŐS-HEILBRONN CONJECTURE

Vsevolod F. Lev

Institute of Mathematics Hebrew University Jerusalem 91904, Israel

seva@math.huji.ac.il

Submitted: June 18, 1998; Accepted: January 29, 2000

ABSTRACT. In 1980, Erdős and Heilbronn posed the problem of estimating (from below) the number of sums a+b where $a \in A$ and $b \in B$ range over given sets $A, B \subseteq \mathbb{Z}/p\mathbb{Z}$ of residues modulo a prime p, so that $a \neq b$. A solution was given in 1994 by Dias da Silva and Hamidoune. In 1995, Alon, Nathanson and Ruzsa developed a polynomial method that allows one to handle restrictions of the type $f(a, b) \neq 0$, where f is a polynomial in two variables over $\mathbb{Z}/p\mathbb{Z}$.

In this paper we consider restricting conditions of general type and investigate groups, distinct from $\mathbb{Z}/p\mathbb{Z}$. In particular, for $A, B \subseteq \mathbb{Z}/p\mathbb{Z}$ and $\mathcal{R} \subseteq A \times B$ of given cardinalities we give a sharp estimate for the number of distinct sums a+b with $(a,b) \notin \mathcal{R}$, and we obtain a partial generalization of this estimate for arbitrary Abelian groups.

1. Background: Mapping restrictions

For two subsets A and B of the set of elements of a group G we write

$$A \dotplus B = \{a + b \colon a \in A, \ b \in B, \ a \neq b\}.$$

(The group $G = \mathbb{Z}/p\mathbb{Z}$ of residues modulo a prime p was historically first to emerge in this context, hence the additive notation.) In other words, $A \dotplus B$ is the set of all elements of G, representable as a sum of two distinct elements from A and B.

The Erdős-Heilbronn conjecture (see [5, p. 95]), resolved (affirmatively) in [4] (cf. also [1, 2]) is the following.

¹⁹⁹¹ Mathematics Subject Classification. Primary: 11B75; Secondary: 11P99, 05C25, 05C35. Partially supported by the Edmund Landau Center for Research in Mathematical Analysis and Related Areas, sponsored by the Minerva Foundation (Germany).

Conjecture 1 (Erdős and Heilbronn). For any two sets $A, B \subseteq \mathbb{Z}/p\mathbb{Z}$,

(1)
$$|A + B| \ge \min\{|A| + |B| - 3, p\}.$$

The set $A \dotplus B$ is obtained from $A + B = \{a + b : a \in A, b \in B\}$ by excluding those sums with b = a. It seems plausible that (1) remains valid even if the sums to be excluded are chosen according to a more general pattern.

Specifically, given a mapping $\tau \colon A \to B$, we define A + B to be the set of all the sums a + b such that $b \neq \tau(a)$:

$$A + B = \{a + b : a \in A, b \in B, b \neq \tau(a)\}.$$

In [3], we conjectured the following.

Conjecture 2 (Lev). Let A and B be subsets of $\mathbb{Z}/p\mathbb{Z}$ satisfying $|A| \leq |B|$, and let $\tau: A \to B$ be an arbitrary mapping from A to B. Then

$$|A + B| \ge \min\{|A| + |B| - 3, p\}.$$

It turns out, however, that this latter conjecture was too optimistic. Fix two integers $k, d \ge 1$ such that $(2k+1)(2d-1) \le 2p+1$ and let

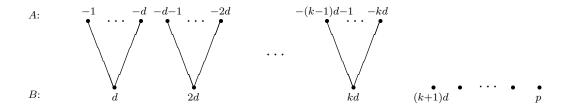
$$A = \{-1, -2, -3, \dots, -kd\} \pmod{p},$$

$$B = \{0, d, 2d, \dots, (k+1)d, (k+1)d + 1, \dots, p-1\} \pmod{p}.$$

Furthermore, define $\tau \colon A \to B$ by

$$\tau(-td+r) = td; \quad 1 < t < k, \ 0 < r < d.$$

This construction can be illustrated by the diagram below:



We have then

$$|A| + |B| = kd + k + (p - (k+1)d + 1) = p + k + 1 - d,$$

while

$$A + B = \{d, d+1, d+2, \dots, p-1\},\$$

and therefore

$$|A + B| = p - d = |A| + |B| - k - 1.$$

This shows that some additional conditions are necessary in order for $|A + B| \ge |A| + |B| - 3$ to hold.

The situation might be somewhat better if τ is *injective*. In this case, we were only able to construct A, B and τ so that |A| + |B| is as large as $\lfloor 4p/3 \rfloor$, and yet |A + B| = p - 2. Specifically, assume for definiteness p = 3m + 2 (the case $p \equiv 1 \pmod{3}$) can be dealt with similarly) and let

$$A = \{1, 3, 4, 6, 7, \dots, 3m, 3m + 1\} \pmod{p},$$

$$B = \{0, -1, -3, -4, -6, -7, \dots, -3m + 2, -3m\} \pmod{p}.$$

Furthermore, define $\tau(a_i) = b_i$, where a_i and b_i are *i*th elements of A and B, respectively, in the above indicated order. Then

$$|A| + |B| = 2(2m+1) = \frac{1}{3}(4p-2) = \lfloor 4p/3 \rfloor,$$

while $A \stackrel{\tau}{+} B$ consists of all residues modulo p, except 1 and 2, and therefore

$$|A + B| = p - 2.$$

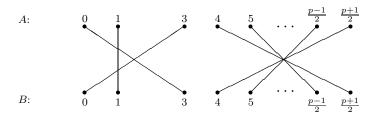
The following example (based on a suggestion of A. Dudek for p = 11) gives the largest known value of |A| + |B| subject to |A + B| = p - 3. Let

$$A = B = \{0, 1, 3, 4, 5, \dots, (p-1)/2, (p+1)/2\},$$

$$\tau(0) = 3, \ \tau(1) = 1, \ \tau(3) = 0,$$

$$\tau(a) = (p+9)/2 - a \text{ for } a = 4, 5, \dots, (p-1)/2, (p+1)2.$$

The corresponding diagram:



Here |A| + |B| = p + 1, and as indicated above, |A + B| = p - 3 (notice that 2, 3 and (p+9)/2 are not in A + B).

Is it true that

(2)
$$|A + B| \ge \begin{cases} |A| + |B| - 3, & \text{if } |A| + |B| \le p, \\ p - 3, & \text{if } |A| + |B| = p + 1, \\ p - 2, & \text{if } |A| + |B| \ge p + 2 \end{cases}$$

for any injective τ ? Though this may not be the case, there is little doubt that it is close to the truth. Finding the best possible estimate of this kind seems to require fresh ideas and is of certain interest.

2. Summary of results

Below, we discuss the results that will be proved in Section 3.

We first bring into consideration restricting conditions of a more general sort. Let $A, B \subseteq G$ be subsets of a group G, and let $\mathcal{R} \subseteq A \times B$ be any subset of the Cartesian product $A \times B$. We define A + B to be the set of all sums a + b, such that $(a, b) \notin \mathcal{R}$:

$$A + B = \{a + b \colon a \in A, b \in B, (a, b) \notin \mathcal{R}\}.$$

To simplify the notation, we write throughout the rest of the paper

$$m = |A|, n = |B|, r = |\mathcal{R}|,$$

and we tacitly assume r > 0 (that is, $\mathcal{R} \neq \emptyset$). Plainly, when \mathcal{R} is induced by a mapping $\tau \colon A \to B$, we have r = m.

Our main result for the group $\mathbb{Z}/p\mathbb{Z}$ is the following.

Theorem 1. Let $A, B \subseteq \mathbb{Z}/p\mathbb{Z}$, and let $\mathcal{R} \subseteq A \times B$. Assume for definiteness $m \leq n$. Then

$$|A \overset{\mathcal{R}}{+} B| \ge \begin{cases} m + n - 2\sqrt{r} - 1, & \text{if} \quad m + n \le p + \sqrt{r} \text{ and } \sqrt{r} \le m, \\ p - \frac{r}{m + n - p}, & \text{if} \quad m + n \ge p + \sqrt{r}, \\ n - \frac{r}{m}, & \text{if} \quad \sqrt{r} \ge m. \end{cases}$$

Observe, that $m+n \ge p+\sqrt{r}$ and $\sqrt{r} \ge m$ can occur simultaneously only when n=p and $m=\sqrt{r}$, in which case the two last estimates of Theorem 1 coincide.

Theorem 1 is extremely sharp and in fact, establishes the minimum possible value of the cardinality of the restricted sum A + B. To see this, consider the following example.

Example 1. Let $A = \{0, \ldots, m-1\} \pmod{p}$ and $B = \{0, \ldots, n-1\} \pmod{p}$, where $1 \le m \le n \le p$. Fix a positive integer k such that

$$\frac{m+n-1-p}{2} \le k \le \frac{m+n-1}{2}$$

and define

$$\mathcal{R} := \{(a,b) \colon a \in A, b \in B, \ a+b \notin [k, m+n-2-k] \pmod{p}\}.$$

(Notice, that \mathcal{R} "eliminates" sums a+b with minimal number of representations. The condition that k is positive ensures that $\mathcal{R} \neq \emptyset$.) We have then $A+B \subseteq [k, m+n-2-k]$ (mod p), whence

$$|A + B| \le m + n - 2k - 1.$$

Now, if k is chosen to satisfy $m \le k \le (m+n-1)/2$, one can verify that

$$r = m(2k + 1 - m) > m^2,$$

and by (3),

$$|A + B| \le n - \frac{r}{m}.$$

If $m+n-p \le k \le m-1$, then

$$r = k(k+1) < m^2$$
, $m+n < p+k < p+\sqrt{r}$;

by (3),

$$|A + B| < m + n - 2\sqrt{r}.$$

Finally, if $(m+n-p-1)/2 \le k \le m+n-p-1$, then

$$r = (p + 2k + 1 - m - n)(m + n - p) < (m + n - p)^{2};$$

it follows that $m+n>p+\sqrt{r}$ and

$$|A \overset{\mathcal{R}}{+} B| \le p - \frac{r}{m+n-p}.$$

We now turn to generalizations onto groups, distinct from $\mathbb{Z}/p\mathbb{Z}$. The second estimate of Theorem 1 has an analog even in the non-commutative case.

Theorem 2. Let $A, B \subseteq G$ be subsets of a finite group G of order q = |G|, and let $\mathcal{R} \subseteq A \times B$. Suppose that $m + n \ge q + 1$. Then

$$|A \overset{\mathcal{R}}{+} B| \ge q - \frac{r}{m + n - q}.$$

Corollary 1. Let G, A, B and \mathcal{R} be as in Theorem 2. Assume that $m + n \ge (1 + \varepsilon)q$ and $r < C \min\{m, n\}$ for some $\varepsilon, C > 0$.

$$|A + B| \ge q - C\varepsilon^{-1}$$

Proof.

$$\frac{r}{m+n-q} \leq \frac{C}{2} \, \frac{m+n}{(m+n)-q} \leq \frac{C}{2} \, \frac{1}{1-(1+\varepsilon)^{-1}} = \frac{C}{2} \left(1+\frac{1}{\varepsilon}\right) < C\varepsilon^{-1}.$$

A refinement is possible when \mathcal{R} is induced by an injective mapping and the sum m+n only slightly exceeds q.

Theorem 3. Let $A, B \subseteq G$ be subsets of a finite group G of order q = |G|, and let $\tau \colon A \to B$ be an injective mapping from A to B. Suppose that $m + n \ge q + 1$. Then

$$|A + B| > q - \sqrt{q} - \frac{1}{2}.$$

In a certain (rather narrow) range of m, n and in the particular case of \mathcal{R} induced by an injective mapping, Theorem 3 improves Theorem 1.

To deal with the generalization of the most important and most difficult case of Theorem 1 — that of small m+n — we make a simplifying assumption that \mathcal{R} satisfies the following conditions:

- (a) to any fixed $a_0 \in A$ there corresponds at most one $b \in B$ such that $(a_0, b) \in \mathcal{R}$;
- (b) to any fixed $b_0 \in B$ there corresponds at most one $a \in A$ such that $(a, b_0) \in \mathcal{R}$.

We note that these conditions automatically hold when \mathcal{R} is induced by an injective mapping; in general, they are not vital, but make possible certain simplifications. Furthermore, for real L > 0 we consider an additional condition:

(c) for any $c \in G$ there are at most L pairs $(a, b) \in \mathcal{R}$ such that a + b = c.

The relevance of this condition for the estimates of |A + B| is hinted to by [6, Conjecture 2] and [6, Theorem 3]: when \mathcal{R} is induced by the equality relation, L can be chosen to be the "doubling constant" of [6].

Our next result is parallel to [6, Theorem 3]. The difference is that in [6] we were only concerned with the "classical" restriction $b \neq a$ and considered only the case B = A; on the other hand, the latter allowed us to cover non-commutative groups.

Theorem 4. Let G be an Abelian group, let $A, B \subseteq G$ be subsets of G, and let \mathcal{R} satisfy conditions (a)–(c). Suppose that $A + B \neq A + B$. Then

(4)
$$|A + B| > (1 - \delta)(m + n) - (L + 2),$$

where

$$\delta = \frac{mn}{(m+n)^2} \le \frac{1}{4}.$$

The condition $A + B \neq A + B$ may look odd at first sight and is worth explanation. The point is that there is such a powerful tool as Kneser's theorem to estimate the number of elements of the *non-restricted* sum A + B from below. If A + B = A + B, this theorem automatically yields lower-bound estimates for the number of elements of the restricted sum A + B; Theorem 4 deals with the complementary case $A + B \neq A + B$. To be more specific, if (4) fails, while A + B = A + B, then it follows immediately from Kneser's theorem that A and B posses a very rigid structure:

– either there exist elements $a \in A$, $b \in B$ and a subgroup $H \subseteq G$ such that $A \subseteq a + H$, $B \subseteq b + H$, $m + n \ge 4(|H| + L + 2)/3$, and

$$A + B = A + B = a + b + H;$$

- or there exist elements $a_1, a_2 \in A$, $b \in B$ and a subgroup $H \subseteq G$ such that $A \subseteq (a_1 + H) \cup (a_2 + H)$, $B \subseteq b + H$, $m + n \ge 8(|H| + L + 2)/3$, and

$$A + B = A + B = (a_1 + b + H) \cup (a_2 + b + H),$$

- or there exist elements $a \in A$, $b_1, b_2 \in B$ and a subgroup $H \subseteq G$ such that $A \subseteq a + H$, $B \subseteq (b_1 + H) \cup (b_2 + H)$, $m + n \ge 8(|H| + L + 2)/3$, and

$$A + B = A + B = (a + b_1 + H) \cup (a + b_2 + H).$$

The coefficient $1 - \delta$ in Theorem 4 can be slightly improved using the methods of [6]; in particular, for B = A it can be increased to $(\sqrt{5} + 1)/4 \approx 0.80$.

The proofs of Theorems 1–4 are mostly combinatorial, with a somewhat surprising interference of graph theory in the proof of Theorem 3 — see also Section 4, the Conclusion.

3. Proofs

Proof of Theorem 1. For i = 1, 2, ... we denote by N_i the number of residues $c \in A + B$ with at least i representations of the form c = a + b ($a \in A$, $b \in B$), and by N'_i the number of residues $c \in (A + B) \setminus (A + B)$ with at least i representations of this form. Obviously, $N_i - N'_i$ counts the number of elements of A + B with at least i representations, whence $N_i - N'_i \leq |A + B|$ and

(5)
$$t|A + B| \ge (N_1 - N_1') + \dots + (N_t - N_t')$$

for any integer $t \geq 1$.

Now, by Pollard's theorem (see [7]) we have

$$N_1 + \cdots + N_t \ge t \min\{p, m+n-t\},$$

provided $t \leq m$, and at the same time, clearly

$$N'_1 + \dots + N'_t \le N'_1 + \dots + N'_t + \dots = \sum_{c \in (A+B) \setminus (A + B)} \nu(c) \le r,$$

where $\nu(c)$ is the number of representations of c. Comparing to (5) we conclude that

$$t|A + B| \ge t \min\{p, m + n - t\} - r,$$

 $|A + B| \ge \min\{p - r/t, m + n - (t + r/t)\},$

and it remains to optimize in t by choosing

$$t = \begin{cases} \lceil \sqrt{r} \rceil, & \text{if } m+n \le p + \sqrt{r} \text{ and } \sqrt{r} \le m, \\ m+n-p, & \text{if } m+n \ge p + \sqrt{r}, \\ m, & \text{if } \sqrt{r} \ge m. \end{cases}$$

Proof of Theorem 2. Let S be the complement of A + B in G, so that |S| = q - |A + B|. By the Dirichlet boxing principle, to any $s \in S$ there correspond at least m + n - q pairs (a, b) (with $a \in A$, $b \in B$) such that a + b = s, and for any such pair we have $(a, b) \in \mathcal{R}$. Totally, we have at least |S|(m+n-q) pairs $(a, b) \in \mathcal{R}$. On the other hand, the number of these pairs is r, whence $|S| \leq r/(m+n-q)$, and the result follows. \square

Proof of Theorem 3. We define S as above, and we want to prove that $|S| < \sqrt{q} + 1/2$. It is convenient to use graph-theoretic terminology. Consider the |S|-regular bipartite graph Γ on two disjoint copies of G, obtained by joining each vertex $x \in G$ of the first copy to the |S| vertices -x + s; $s \in S$ of the second copy. Formally, we write

$$\Gamma = (X \cup Y, E); \quad E = \{(x, y) : x \in X, y \in Y, x + y \in S\},\$$

where X and Y are thought of as two disjoint copies of G. Furthermore, consider the subgraph $\Gamma_0 \subseteq \Gamma$, induced by all elements of A in the first copy and all elements of B in the second copy:

$$\Gamma_0 = \langle (X \cap A) \cup (Y \cap B) \rangle.$$

We claim that Γ_0 contains no paths of length two. Indeed, a path x_1, y, x_2 with $x_1, x_2 \in X \cap A$ and $y \in Y \cap B$ would mean $x_1 + y \in S$ and $x_2 + y \in S$, which is impossible: either $\tau(x_1) \neq y$ (in which case $x_1 + y \notin S$), or $\tau(x_2) \neq y$ (in which case $x_2 + y \notin S$). Similarly, a path of the type y_1, x, y_2 cannot occur in Γ_0 as $\tau(x) = y_1$ and $\tau(x) = y_2$ cannot happen simultaneously.

Our next observation is that Γ contains no rectangles. Indeed, any single rectangle x_1, y_1, x_2, y_2 can be translated to produce q rectangles

$$x_1 + u, -u + y_1, x_2 + u, -u + y_2; u \in G,$$

each containing at most two vertices of Γ_0 : a subgraph induced by any three vertices of a rectangle necessarily contains a path of length two. Summation over all $u \in G$ gives $2(|A| + |B|) \leq 2q$, contradicting the assumptions.

We now essentially repeat an Erdős' argument to show that if Γ contains no rectangles, then |S| (the degree of Γ) is small. We first count all paths of the form x_1, y, x_2 ($x_1, x_2 \in X, y \in Y$) in Γ . Obviously, there are totally $q\binom{|S|}{2}$ such paths, as any vertex $y \in Y$ participates in $\binom{|S|}{2}$ paths. On the other hand, there are only $\binom{q}{2}$ pairs (x_1, x_2) ; $x_i \in X$ with $x_1 \neq x_2$. Since no two distinct paths can share a common pair (this would yield a rectangle), we have

(6)
$$q {|S| \choose 2} \le {q \choose 2},$$
$$|S|^2 - |S| + 1 \le q$$

whence $|S| < \sqrt{q} + 1/2$, as required.

There is another way to complete the proof by making a funny observation that S is a Sidon set in G: an equality $s_1 - s_2 = s'_1 - s'_2$ with $s_1 \neq s_2$, $s_1 \neq s'_1$ creates a rectangle $s_1, 0, s_2, -s_2 + s'_2 = -s_1 + s'_1$. It is easy to verify, however, that the cardinality of any Sidon set $S \subseteq G$ satisfies (6).

Proof of Theorem 4. We break the proof in three steps.

i) Since $A + B \neq A + B$, there exist $a_0 \in A$ and $b_0 \in B$ such that $c = a_0 + b_0 \notin A + B$. Then

$$|(A-b_0)\cap (a_0-B)|\leq L,$$

as any equality $a - b_0 = a_0 - b$ gives rise to the representation c = a + b with $(a, b) \in \mathcal{R}$ (in view of $c \notin A + B$), and there are at most L such representations. Letting $A - B = \{a - b : a \in A, b \in B\}$ we obtain

(7)
$$|A - B| \ge |(A - b_0) \cup (a_0 - B)| \ge m + n - L.$$

ii) Fix any $c = a_0 - b_0 \in A - B$ (where $a_0 \in A$, $b_0 \in B$) and let

$$A_0 = \{ a \in A \colon (a, b_0) \notin \mathcal{R} \}, \quad B_0 = \{ b \in B \colon (a_0, b) \notin \mathcal{R} \},$$

so that $|A_0| \ge m-1$ and $|B_0| \ge n-1$ by the conditions (a) and (b). Write $\nu(c)$ for the number of representations c = a - b ($a \in A, b \in B$). Then

(8)
$$\nu(c) \ge |(a_0 + B_0) \cap (A_0 + b_0)| \ge |A_0| + |B_0| - |A + B|$$

$$\ge m + n - 2 - |A + B|,$$

as $(a_0 + B_0) \cup (A_0 + b_0) \subseteq A + B$.

iii) By (7) and (8),

$$mn = \sum_{c \in A-B} \nu(c) \ge (m+n-L)(m+n-2-|A \overset{\mathcal{R}}{+} B|),$$
$$|A \overset{\mathcal{R}}{+} B| \ge m+n-2-\frac{mn}{m+n-L}$$
$$> m+n-\frac{mn}{m+n}-(L+2),$$

the latter inequality being equivalent to mn < (m+n)(m+n-L), which follows from $L < (1-\delta)(m+n)$ — otherwise, the assertion of the theorem is trivial.

The result follows.

4. Conclusion

We re-state here explicitly several problems that remain open.

Does (2) hold for any two sets $A, B \subseteq \mathbb{Z}/p\mathbb{Z}$ and any injective mapping $\tau \colon A \to B$? In particular, is it true that $|A + B| \ge p - 2$, provided $|A| + |B| \ge p + 2$? Do similar estimates hold when A and B are subsets of an arbitrary finite group? If some of the answers are negative, what are the best possible estimates of this sort?

A slight modification of the approach used in the proof of Theorem 3 shows that for $|A| + |B| \ge p + 2$, this problem can be reformulated in the graph-theoretic language as follows.

Let $c \neq 0, 1$ be any fixed residue modulo p. Consider a cubic bipartite graph $\Gamma(c)$ on two copies of $\mathbb{Z}/p\mathbb{Z}$ such that any vertex x of the first copy is adjacent to the vertices x, x+1 and x+c of the second copy. Is it true that any induced subgraph of order p+2 of any such $\Gamma(c)$ contains a path of length two?

The answer is certainly positive if c = -1, 2 or (p+1)/2 (in which cases $\Gamma(c)$ contains a rectangle). In general, the situation is not clear, however.

The major open problem for generic restriction is that of improving the coefficient $1-\delta$ in Theorem 4. Quite likely, this coefficient can be replaced by 1 or at least by $1-\varepsilon$ for any positive ε , provided r is sufficiently (in terms of ε) small compared to m+n.

ACKNOWLEDGMENT

The (core of the) present form of Theorem 1 and an idea incorporated in the proof of Theorem 3 are due to Noga Alon, whom we are greatly indebted for this contribution.

References

- [1] N. Alon, M.B. Nathanson and I.Z. Ruzsa, Adding distinct congruence classes modulo a prime, *American Math. Monthly*, **102** (1995), 250–255.
- [2] N. Alon, M.B. Nathanson and I.Z. Ruzsa, The Polynomial Method and Restricted Sums of Congruence Classes, *J. Number Theory*, **56** (1996), 404–417.
- [3] Y. Bilu, V.F. Lev and I.Z. Ruzsa, Rectification principles in additive number theory, *Disc.* and Comput. Geometry, **19** (1998), 343–353.
- [4] J.A. DIAS DA SILVA and Y.O. HAMIDOUNE, Cyclic spaces for Grassmann derivatives and additive theory, *Bull. London Math. Soc.* **26** (1994), 140–146.
- [5] P. Erdős and R.L. Graham, Old and new problems and results in combinatorial number theory, L'Enseignement Mathématique, Geneva, 1980.
- [6] V.F. Lev, Restricted set addition in groups. I. The classical setting, *Journal of the London Mathematical Society*, to appear.
- [7] J.M. Pollard, A generalization of the theorem of Cauchy and Davenport, J. London Math. Soc. 2 (8) (1974), 460–462.