

Hitting k primes

comment, rough draft, remark199

1 Results

Let $S = (d_1, d_2, d_3, \dots)$ be an infinite sequence of rolls of independent fair dice. Thus the d_i are independent, identically distributed random variables, each uniformly distributed on the integers $\{1, 2, \dots, 6\}$. For each $i \geq 1$ put $s_i = \sum_{j=1}^i d_j$. The sequence S hits a positive integer x if there exists an i so that $s_i = x$. In that case it hits x in step i .

For any positive integer k , let $L_k = L_k(S)$ be the random variable whose value is the smallest i so that the sequence S hits k primes during the first i steps (∞ is there is no such i , but it is easy to see that with probability 1 there is such i). The random variable L_1 is introduced and studied in [1], see also [3] for several generalizations.

I think there have also been some additional follow-up papers that may be mentioned (?)

Here we consider the random variable L_k for larger values of k , focusing on the estimate of its expectation.

1.1 Computational results

The value of the expectation of L_k for $k \leq 30$ is given in the following table.

Insert the table here.

The table suggests that the asymptotic value of this expectation is $(1 + o(1))k \log k$, where the $o(1)$ -term tends to zero as k tends to infinity, and the logarithm here and throughout the manuscript is in the natural basis. This is confirmed in the results stated and proved below.

1.2 Asymptotic results

Theorem 1.1. *For any fixed positive reals ε, δ there exists $k_0 = k_0(\varepsilon, \delta)$ so that for all $k > k_0$ the probability that $|L_k - k \log k| > \varepsilon k \log k$ is smaller than δ .*

Theorem 1.2. *For any fixed $\varepsilon > 0$ and any $k > k_0(\varepsilon)$, the expected value of the random variable L_k satisfies $|E(L_k) - k \log k| < \varepsilon k \log k$.*

2 Proofs

Lemma 2.1. *There are fixed positive C and μ so that the following holds. Let $S = (d_1, d_2, \dots)$ be a random sequence as above. For any positive integer x , let $p(x)$ denote the probability that*

S hits x . Then $|p(x) - 2/7| \leq C(1 - \mu)^x$, that is, as x grows, $p(x)$ converges to the constant $2/7$ with an exponential rate.

Proof. Define $p(-5) = p(-4) = p(-3) = p(-2) = p(-1) = 0$, $p(0) = 1$ and note that for every $i \geq 1$,

$$p(i) = \frac{1}{6} \sum_{j=1}^6 p(i-j)$$

Indeed, S hits i if and only if the last number it hits before i is $i-j$ for some $j \in \{1, \dots, 6\}$, and the die rolled after that gives the value j . The probability of this event for each specific value of j is $p(i-j) \cdot (1/6)$, providing the equation above. (Note that the definition of the initial values is consistent with this reasoning, as before any dice rolls the initial sum is 0). Thus, the sequence $(p(i))$ satisfies the homogeneous linear recurrence relation given above. The characteristic polynomial of that is

$$P(z) = z^6 - \frac{1}{6}(z^5 + z^4 + z^3 + z^2 + z + 1).$$

One of the roots of this polynomial is $z = 1$, and its multiplicity is 1 as the derivative of $P(z)$ does not vanish at 1. It is also easy to check that the absolute value of each of the other roots λ_j , $2 \leq j \leq 6$ of $P(z)$ is at most $1 - \mu$ for some absolute positive constant μ . Therefore, there are constants c_j so that

$$p(i) = c_1 \cdot 1^i + \sum_{j=2}^6 c_j \lambda_j^i,$$

implying that

$$|p(i) - c_1| \leq C(1 - \mu)^i$$

for some absolute constant C . It remains to compute the value of c_1 . By the last estimate, for any positive n ,

$$\left| \sum_{i=1}^n p(i) - c_1 n \right| \leq C/(1 - \mu).$$

Note that the sum $\sum_{i=1}^n p(i)$ is the expected number of integers in $[n] = \{1, 2, \dots, n\}$ hit by the sequence S . For large n , this number is clearly $(1 + o(1))(2/7)n$, by standard estimates for distributions of sums of independent bounded random variables, see, e.g., [2], Theorem A.1.16. Dividing by n and taking the limit as n tends to infinity shows that $c_1 = 2/7$, completing the proof. \square

The next simple lemma shows that for any integers $x_1 < x_2 < \dots < x_r$ that are far from each other, the events that the random sequence S hits x_i are nearly independent.

Lemma 2.2. *For any positive integers $x_1 < x_2 < \dots < x_r$, the probability that the random sequence S hits all x_i is exactly*

$$p(x_1)p(x_2 - x_1)p(x_3 - x_2) \cdots p(x_r - x_{r-1}).$$

Therefore, if each difference $(x_i - x_{i-1})$ is at least s , then the probability of this event deviates from $(\frac{2}{7})^r$ by at most $(\frac{2}{7})^{r-1}Cr(1 - \mu)^s$.

Proof. (Sketch) The conditional probability of the sequence to hit x_{i+1} given that it hit already x_1, \dots, x_i is exactly $p(x_{i+1} - x_i)$, since the sequence starting at x_i has the same distribution as S . The desired estimate follows from the assertion of Lemma 2.1. \square

Using the two lemmas above we next show that the number of primes hit by the first f steps of the random sequence S is close to $2/7$ times the number of primes smaller than $3.5f$ with high probability.

Theorem 2.3. *Let $\pi(x)$ denote the number of primes smaller than x , and let $Y(f)$ denote the number of primes hit by the random sequence S during the first f steps. Then, for any fixed $\varepsilon > 0$ and any (large) constant t , and for any $f > f_0(\varepsilon, t)$, the probability that $Y(f)$ deviates from $(2/7)\pi(3.5f)$ by more than $\varepsilon(2/7)\pi(3.5f)$ is smaller than f^{-t} .*

Proof. (rough sketch) For large f , with probability larger than $1 - f^{-t}$ the sum $\sum_{i=1}^f d_i$ deviates from its expectation $3.5f$ by less than $(\varepsilon/3)3.5f$. Split all the primes smaller than $(1 - \varepsilon/3)3.5f$ into, say, \sqrt{f} groups of nearly equal sizes, where the difference between any two elements in the same group is at least \sqrt{f} . Using Lemma 2.2 it follows that for each fixed group of size g , the number of primes of the group hit by S is within an ε -fraction of its expectation. This is done by computing, say, the first $4t$ moments of this random variable, observing that these are very close to the same moments of a random variable which is the sum of g independent indicator random variables, each being 1 with probability $2/7$. By considering the expectation of the $4t$ -th power of the difference between this random variable and its expectation, this implies the desired concentration within each group, and the triangle-inequality supplies the required estimate for the union of all groups. The contribution of the primes between $(1 - \varepsilon/3)3.5f$ and $(1 + \varepsilon/3)3.5f$ is small, by the known results about the distribution of primes, and the contribution of the primes larger than $(1 + \varepsilon/3)3.5f$ to the expectation is negligible, since the probability that the sequence reaches these numbers within the first f steps is tiny. This implies the assertion of the Theorem. \square

The assertions of Theorem 1.1 and Theorem 1.2 can be easily deduced from that of Theorem 2.3.

3 Concluding remarks and extensions

- **Extensions for biased r -sided dice and arbitrary subsets of the integers.** The proofs in the previous section use very little of the specific properties of the primes and the specific distribution of each d_i . It is easy to extend the result to any r -sided dice with an arbitrary discrete distribution on $[r]$ in which the values obtained with positive probabilities do not have any nontrivial common divisor. The constants 3.5 and $2/7$ will then have to be replaced by the expectation of the random variable d_i and by its reciprocal, respectively. Similarly, we can replace the set of integers in which we count hits by an arbitrary set of positive integers, as long as its distribution satisfies some mild smoothness assumptions. We omit the details.
- **Heuristic suggestion for a more precise expression for $E(L_k)$.** We may state here a possible more precise expression for $E(L_k)$, which may be close to $k(\log k + \log \log k + c_1) + c_2$. This is justified very loosely by the heuristic argument described here together with the behavior of $\pi(n)$, and is also roughly consistent with the experimental evidence. Better to state it only as a possible guess and mainly raise the question of finding a more accurate estimate for the error term in the expectation.
- Can add the definition and conjecture that there are infinitely many AMMZ integers (after deciding if this is indeed the name we want to suggest). The known results and conjectures about prime gaps (specifically Cramér's Conjecture) and the fact that the difference between $E(L_{k+1})$ and $E(L_k)$ is close to $\log k$ (at least by our heuristics, supported by the computation), suggest that the function $2\pi(\text{trunc}(\frac{7}{2} \cdot L[i]))$ will stay an even constant integer for some $(2/7) \log(L[i]) = (1 + o(1))(2/7) \log i$ consecutive values of i around i for infinitely many values of i . As the denominator $7i$ will go through some $(2/7) \log i$ consecutive even values for this range, there is, possibly, a non-negligible chance that one of these will be equal to the numerator. Of course this depends on quite a few heuristic conjectures, including Cramér's. The computational results in the table we have do not handle sufficiently large i to test if this heuristics is valid. In particular, I think that in the range in the table there are no two consecutive values in which the numerator in our formula is the same - this should happen later infinitely often.

References

- [1] Noga Alon and Yaakov Malinovsky, Hitting a prime in 2.43 dice rolls (on average) The American Statistician, Volume 77 Issue 3 (2023), 301-303.

- [2] Noga Alon and Joel H. Spencer, *The Probabilistic Method, Fourth Edition*, Wiley, 2016, xiv+375 pp.
- [3] Lucy Martinez and Doron Zeilberger, How many dice rolls would it take to reach your favorite kind of number? arXiv:2302.00143, 2023.