

LATTICE POINTS ON AND NEAR CIRCLES

CRAIG CHEN

ABSTRACT. This essay studies the problem of counting the number of integer lattice points on and within a neighborhood of circles and ellipses. Of particular interest is the problem of counting how many lattice points can lie on an arc of length R^θ of the circle centered at the origin with radius R for $\theta \in [0, 1)$. While this problem is interesting in its own right, solutions to these types problems have consequences in the field of PDEs. One such example is the local wellposedness of the initial value problem for certain non-linear Schrödinger equations on general Tori.

1. INTRODUCTION

Something that has engaged mathematicians for centuries is Gauss's eponymous question that asks: how many lattice points are within a circle centered at the origin of radius R ? In the past century, significant progress has been made towards answering this question. There are explicit formulas for this number; however, since the number of lattice points in the circle is approximately the area πR^2 , the problem has taken the form of trying to precisely bound the error of this approximation, and this has proven to be very difficult. We are interested in a similar question that is concerned with the lattice points *on*, not within, the circle centered at the origin of radius R . In fact, there is also an explicit formula for the total number of lattice points on a circle of radius R ; however, if we restrict ourselves to an arc of length R^θ , the number of lattice points on such an arc is much less understood. Another interesting question is what happens if we allow a small neighborhood of the circle (or more generally an ellipse).

In this essay, we provide a rigorous proof of the current best bound on the number of lattice points on a small arc of the circle. Throughout the proof, we will discuss the potential reasoning behind each step as well as comment on potential avenues of improvement. We also provide more detailed proofs of two other results bounding the number of lattice points within a neighborhood of an ellipse and of a convex curve (the types of neighborhoods considered in these two results are very different). We assume that the reader is familiar with basic results about prime numbers and complex analysis; for example, Euler's formula $e^{it} = \cos(t) + i \sin(t)$ which yields a beautiful geometric interpretation of exponentiation. We will frequently jump between the geometric and complex-exponential interpretations of lattice points – comfort with these interpretations is probably necessary for getting the most out of the content of this essay.

For fun, we also briefly and very informally go through a connection between our problem and the well-posedness theory of the Schrödinger equation.

Notation: We say that $A \ll_p B$ if $A \leq C_p B$ for some constant C_p that depends only on p . If there is no subscript, then this constant does not depend on any of the variables in the equation. We say that $f = \mathcal{O}(g)$ if $\exists x_0 \in \mathbb{R}$ such that $\forall x \geq x_0$, $|f(x)| \ll g(x)$.

We will use $\|\cdot\|$ to represent multiple functions: in the proof of Theorem 3.1 we exclusively use this to represent the distance to the nearest integer, and later on we will use

this to denote the Euclidean norm on \mathbb{R}^k . In the proofs we will make clear what definition we are using. $B(a, r)$ denotes an open ball centered at a of radius r .

Sometimes, to save space, we will use $n(m)$ as shorthand for $n \pmod{m}$.

2. PRELIMINARIES AND BACKGROUND

This topic of study is inspired by two results of Czech mathematician Vojtech Jarník.

2.1. Jarník's Results.

Theorem 2.1 (Jarník, [Jar26]).

Any convex curve of length L in \mathbb{R}^2 passes through at most

$$\frac{3}{\sqrt[3]{2\pi}}L^{2/3} + \mathcal{O}(L^{1/3})$$

integer lattice points.

Theorem 2.2.

On a circle centered at the origin of radius R , any arc of length $R^{1/3}$ contains at most 2 lattice points.

The proof of Theorem 2.2 is very cute and the high level strategy of this proof is shared by the proof of Theorem 1 of Cilleruelo and Cordoba (our Theorem 3.1) and the proof of a lemma from Bourgain (our Lemma 5.1).

Proof. Suppose there are 3 lattice points on an arc of the circle. If so, we can connect the points to form a triangle. Since each of the lattice points have integer coordinates, the minimum area A of any possible triangle generated by these 3 points is $\frac{1}{2}$. This gives the lower bound

$$A \geq \frac{1}{2}.$$

We now want to upper bound A in terms of R . The next ingredient in the proof is the formula for the circumradius of a triangle. For a triangle with side-lengths a, b, c and area Δ , the circumradius $r = abc/4\Delta$. By assumption, all of the vertices of our triangle lie on a circle of radius R , so we can rearrange this formula to get an expression for A in terms of R .

$$A = \frac{abc}{4R}.$$

If we can upper bound the side-lengths of the triangle in terms of the length of the arc that contains these three lattice points, we'll be finished with the proof. Let $\lambda_1, \lambda_2, \lambda_3$ denote the three lattice points, encountered in that order when traversing the circle CCW. We can upper bound the side-lengths of the triangle by the arc-lengths to get

$$abc < \text{Arc}(\lambda_1, \lambda_2)\text{Arc}(\lambda_2, \lambda_3)\text{Arc}(\lambda_3, \lambda_1) \leq \frac{1}{4}\text{Arc}(\lambda_3, \lambda_1)^3,$$

where the final inequality comes from the fact that $\text{Arc}(\lambda_1, \lambda_2) + \text{Arc}(\lambda_2, \lambda_3) = \text{Arc}(\lambda_1, \lambda_3)$ and the quantity $x(1-x)$ is maximized when $x = 1/2$ (for $x \in [0, 1]$). Combining this bound with the previous bound, we get

$$\frac{1}{2} \leq A \leq \frac{\text{Arc}(\lambda_1, \lambda_3)^3}{16R} \implies \text{Arc}(\lambda_1, \lambda_3) \geq (8R)^{1/3}.$$

Thus, we conclude that any arc of length less than $(8R)^{1/3}$ contains at-most 2 lattice points. \square

To recap, the strategy behind the proof was to first assume that there are $m + 1$ lattice points on some arc, and then to find a minimum length for that arc. Then, anything smaller must have fewer lattice points.

Remark 2.3. Theorem 2.2 can be improved to $(16R)^{1/3}$.

Proof. Since the 3 lattice points are on a circle centered at the origin, it is impossible for the triangle they generate to have area $1/2$. Let (x_j, y_j) denote the coordinates of λ_j . Since all three lattice points $\lambda_1, \lambda_2, \lambda_3$ lie on the same circle, we know that

$$x_1^2 + y_1^2 = x_2^2 + y_2^2 = x_3^2 + y_3^2 = R^2.$$

We now address two cases:

- $R \equiv 0 \pmod{2}$

This implies that $x_j \equiv y_j \pmod{2}$ for $j = 1, 2, 3$. Since there are only *two* possible values for any $n \pmod{2}$, by pidgeonhole principle, there must be $j \neq j' \in \{1, 2, 3\}$ such that $x_j \equiv x_{j'} \pmod{2}$. In other words, there are two lattice points such that the coordinates of the lattice points are of the same parity.

- $R \equiv 1 \pmod{2}$

The argument for this case is essentially the exact same as the argument above. In this case, we get that $x_j \not\equiv y_j \pmod{2}$ for all j . However, by pidgeonhole, we still get that at least two lattice points must have coordinates that are the same parity.

The above cases imply that the average of these two lattice points is another lattice point. Consequently, the triangle between $\lambda_1, \lambda_2, \lambda_3$ is actually the disjoint union of two integer-lattice-point triangles, so the minimum area is actually 1, not $1/2$. \square

3. LATTICE POINTS ON SHORT ARCS OF A CIRCLE

In this section, we provide a rigorous treatment of Theorem 1 from “Trigonometric Polynomials and Lattice Points” by Cilleruelo and Córdoba.

Theorem 3.1 (Cilleruelo & Córdoba, [CC92]).

On any circle centered at the origin with radius R , there are at most k lattice points on any arc of length

$$\sqrt{2}R^{\frac{1}{2} - \frac{1}{4\lfloor k/2 \rfloor + 2}}.$$

3.1. Representation of Lattice Points as Gaussian Integers.

The first step of the proof of Theorem 3.1 is to construct a correspondence between lattice points of a circle and the Gaussian integers $\mathbb{Z}[i]$. When passing over to the complex numbers, one tends to uncover more interesting results; a famous example of this is the Zeta function. Euler initially defined the function for real-valued inputs and it was only when Riemann extended the definition to complex inputs that we were able to appreciate the much deeper consequences of the properties of the Zeta function.

For a given $n \in \mathbb{N}$, we are interested in whether or not we can write n as the sum of two squares since these representations correspond to the lattice points on the circle of radius \sqrt{n} . Notice that this is not possible for all n and that there may be more than one possible representation.

Lemma 3.2. Define the function $r(n)$ as the number of representations of n as the sum of two squares. By the Fundamental Theorem of Arithmetic, n can be uniquely written

$$(3.1) \quad n = 2^\nu \prod_{p_j \equiv 1 \pmod{4}} p_j^{\alpha_j} \prod_{q_j \equiv 3 \pmod{4}} q_j^{\beta_j}.$$

If $\beta_j \equiv 0 \pmod{2}$ for all j , then $r(n) = 4 \prod (1 + \alpha_j)$.

We will not provide a proof of this lemma since it is not essential to our result; a proof of this result can probably be found in most introductory number theory textbooks or online.

Now, for any positive integer n that can be written as the sum of two squares, trivial decompositions included, we can associate n with the gaussian integers z such that $|z|^2 = n$ where $|\cdot|$ denotes the standard complex norm. If the decomposition into two squares is unique $n = a^2 + b^2$, we can compactly represent all the gaussian integers of norm n as $\sqrt{n}e^{2\pi i(\pm\Phi+t/4)}$ where we define $\Phi = \arg(a + bi)/2\pi$ (with $a + bi$ as the point in the first octant) and $t \in \{0, 1, 2, 3\}$. Something interesting to note is the fact that $\Phi < 1/8$ since we are choosing the points in the first octant. The $t/4$ terms corresponds to multiplication of the gaussian integer $a + bi$ by a unit of $\mathbb{Z}[i]$, namely $1, i, -1, -i$. For convenience, we will omit this $t/4$ term in the enumeration that follows since this term is not influenced by the factorization of n .

For general n there will not be a unique decomposition into the sum of squares. As a stepping stone towards the statement for general n , we first provide more explicit representations for primes and prime-powers using Lemma 3.2

(1) $n = 2^\nu$

The gaussian integers with squared-norm 2^ν are

$$2^{\nu/2} e^{2\pi i(\Phi_0)} \quad \text{where } \Phi_0 = \begin{cases} 0 & \nu \equiv 0 \pmod{2} \\ 1/8 & \nu \equiv 1 \pmod{2} \end{cases}.$$

If ν is even, the four gaussian integers are those on real and imaginary axes. For ν odd, these are the points on the lines $y = x$ and $y = -x$ that intersect the circle of radius \sqrt{n} .

(2) $n = q^\beta, q \equiv 3 \pmod{4}$

If β is odd, there are no possible representations as the sum of two squares. Thus, it suffices to consider the case of $q^{2\beta}$. In this case, the number of interest is itself a perfect square, which gives

$$q^\beta e^{2\pi i \cdot 0}.$$

In other words, $\Phi = 0$.

(3) $n = p^\alpha, p \equiv 1 \pmod{4}$

If $p = a^2 + b^2$, we have the following representation $\sqrt{p}e^{2\pi i\Phi} =: \omega$ which is unique since we are choosing the point in the first octant. To generate the presentations of p^α , we can simply multiply together α gaussian integers with squared-norm p . Ignoring multiplication by units for now, there are two elements we can work with, ω and $\bar{\omega}$.

When multiplying α copies of ω together, we may choose k of them to be the conjugate $\bar{\omega}$ instead since this does not change the norm of the product. For each choice of k , we get a different gaussian integer with squared-norm p^α . Up to multiplication by units, the elements of squared-norm p^α can be written

$$\begin{aligned} \prod_{m=1}^k \bar{\omega} \prod_{n=1}^{\alpha-k} \omega &= \prod_{m=1}^k \sqrt{p} e^{2\pi i(-\Phi)} \prod_{n=1}^{\alpha-k} \sqrt{p} e^{2\pi i\Phi} \\ &= p^{\alpha/2} e^{2\pi i((\alpha-k)\Phi - k\Phi)} \\ &= p^{\alpha/2} e^{2\pi i(\alpha-2k)\Phi}. \end{aligned}$$

Since k can range from 0 to α , we get that

$$(\alpha - 2k) \in \{\gamma \in \mathbb{Z} \mid |\gamma| \leq \alpha, \gamma \equiv \alpha \pmod{2}\} =: \Lambda_\alpha.$$

One may wonder if the lattice points enumerated by the expressions above are all of the lattice points within some quadrant of the circle; unfortunately, this is not the case.

Example 3.3. Consider the circle $x^2 + y^2 = 5^2$. We know that 5 splits into $(2+i)(2-i)$ in $\mathbb{Z}[i]$. Given our convention of taking the point in the first octant, we have that $\omega = 2 + i$ in this case. We can then compute that the 3 = (1 + 2) points given by the above representations. These points are:

$$\omega^2 = (3 + 4i), \quad \omega\bar{\omega} = 5, \quad \bar{\omega}^2 = 3 - 4i.$$

Notice that the lattice points $(4 + 3i)$ and $(4 - 3i)$ actually lie *in between* the lattice points listed above! In other words, the arcs spanned by these collections of lattice points may overlap — to be clear, the term collection is referring to the set of lattice points parameterized by Λ_α , and there are four collections that correspond to multiplication by $1, i, -1, -i$.

In general, for

$$n = 2^\nu \prod_{p_j \equiv 1(4)} p_j^{\alpha_j} \prod_{q_j \equiv 3(4)} q_j^{2\beta_j},$$

by multiplying the expressions above, the gaussian integers of squared-norm n can be written as

$$(3.2) \quad \sqrt{n} e^{2\pi i[\Phi_0 + \sum \gamma_j \Phi_j + t/4]}$$

where

$$\gamma_j \in \Lambda_{\alpha_j}, \quad t \in \{0, 1, 2, 3\}, \quad \Phi_0 = \begin{cases} 0 & \nu \text{ is even} \\ 1/8 & \nu \text{ is odd} \end{cases}.$$

Notice that the above implies that there are $\#(\{0, 1, 2, 3\}) \prod \#(\Lambda_{\alpha_j}) = 4 \prod (1 + \alpha_j)$ representations of n as the sum of two squares which agrees with Lemma 3.2.

Proposition 3.4. *The angles Φ_j are linearly independent over \mathbb{Q} .*

Proof. Suppose not, then there exists $\alpha_0, \alpha_1, \dots, \alpha_n \in \mathbb{Q} \setminus \{0\}$ such that

$$\sum_{j=1}^n \alpha_j \Phi_j + \alpha_0 = 0.$$

Without loss of generality, we can consider integer coefficients since we can multiply through by a common denominator. For each Φ_j , let a_j, b_j be the real and imaginary components of the corresponding gaussian integers; since these are the gaussian primes

that divide the integer primes that are 1 mod 4, we know that $a_j, b_j \neq 0$. We then compute

$$\begin{aligned}
\prod_{j=1}^n (a_j + ib_j)^{2\alpha_j} &= \prod_{j=1}^n (\sqrt{p_j} e^{2\pi i \Phi_j})^{2\alpha_j} \\
&= \exp\left(2\pi i 2 \sum_{j=1}^n \alpha_j \Phi_j\right) \cdot \prod_{j=1}^n (\sqrt{p_j})^{2\alpha_j} \\
&= \exp(2\pi i(-2\alpha_0)) \cdot \prod_{j=1}^n (p_j)^{\alpha_j} \\
&= \prod_{j=1}^n (p_j)^{\alpha_j} \\
&= \prod_{j=1}^n ((a_j + ib_j)(a_j - ib_j))^{\alpha_j}
\end{aligned}$$

However, this is a contradiction since we know that $\mathbb{Z}[i]$ is a unique factorization domain. \square

3.2. Bounding the gap between lattice points.

We are now ready to move on to the main part of the proof of Theorem 3.1. This step of the proof proceeds by first supposing that there are $m + 1$ lattice points on an arc of length $\sqrt{2}R^\theta$ where $R = \sqrt{n}$ and

$$n = 2^\nu \prod_{p_j \equiv 1(4)} p_j^{\alpha_j} \prod_{q_k \equiv 3(4)} q_k^{2\beta_k}.$$

Notice that in Equation 3.2, the argument of the representation does *not* depend on the factors that are 3 (mod 4). Furthermore, the factors of 2 influence the lattice points by rotating them all by $\pi/4$ which does not affect the number of points on an arc of a given length. Thus, we can assume without loss of generality that

$$n = \prod_{p_j \equiv 1(4)} p_j^{\alpha_j}.$$

By Equation 3.2, we can describe each of the $m + 1$ lattice points by a set of $\gamma_j \in \Lambda_{\alpha_j}$ for each $p_j \equiv 1 \pmod{4}$ and a value for $t \in \{0, 1, 2, 3\}$. For $s = 1, 2, \dots, m + 1$, we will denote the parameters as $\gamma_{j,s}$ and t_s . For any two distinct lattice points λ_s and $\lambda_{s'}$, $s \neq s'$, we can define the quantity

$$\begin{aligned}
\Psi_{s,s'} &= \sum_j (\gamma_{j,s} - \gamma_{j,s'}) \Phi_j + \frac{t_s - t_{s'}}{4} \\
&= 2 \left(\sum_j \frac{\gamma_{j,s} - \gamma_{j,s'}}{2} \Phi_j + \frac{t_s - t_{s'}}{8} \right)
\end{aligned}$$

This quantity corresponds to the angle between the lattice points λ_s and $\lambda_{s'}$

$$\frac{\lambda_s}{\lambda_{s'}} = e^{2\pi i [\sum \gamma_{j,s} \Phi_j + t_s/4 - (\sum \gamma_{j,s'} \Phi_j + t_{s'}/4)]} = e^{2\pi i \Psi_{s,s'}}.$$

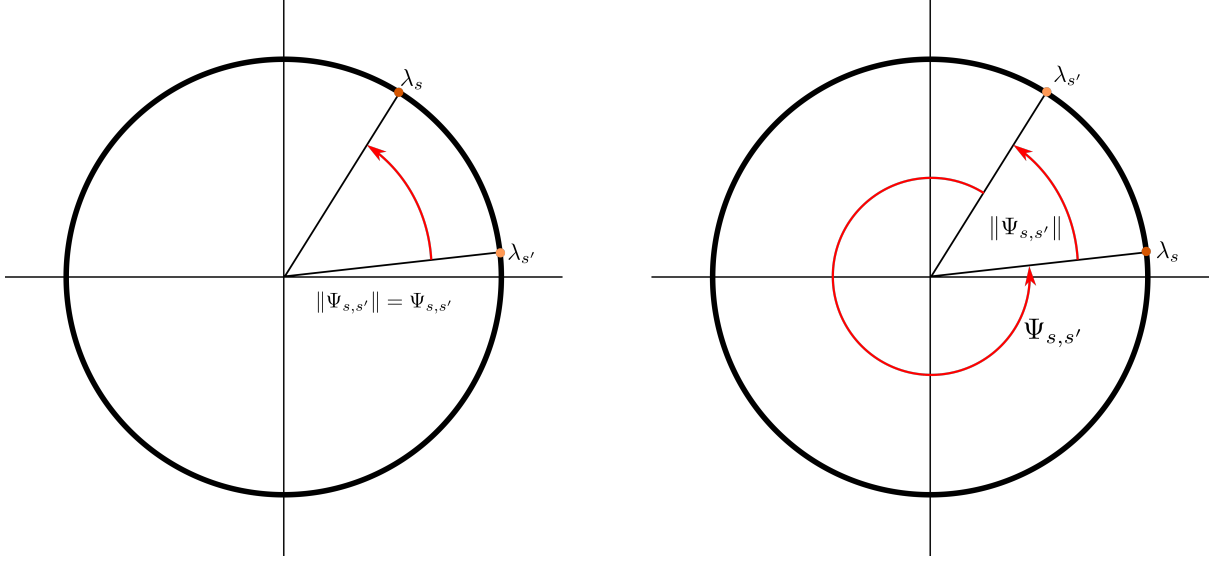


FIGURE 1. $2\pi\|\Psi_{s,s'}\|$ corresponds to the central angle between the lattice points.

Now, define $\|\cdot\| : \mathbb{R} \rightarrow [0, 1/2]$ as the distance to the nearest integer. Then, the quantity $2\pi\|\Psi_{s,s'}\|$ corresponds to the *central* angle between the lattice points λ_s and $\lambda_{s'}$ (see Figure 1). It is worth noting that $\|x/2\| \neq \|x\|/2$ in general, for example take $x = 3/4$.

Consider the quantity $\Psi_{s,s'}/2$. By Proposition 3.4, we know that if $\Psi_{s,s'}/2$ is an integer, then $t_s = t_{s'}$ and $\gamma_{j,s} = \gamma_{j,s'}$ for all j . This implies that $\lambda_s = \lambda_{s'}$. Notice that $(\gamma_{j,s} - \gamma_{j,s'})/2$ is always an integer.

- If $t_s \equiv t_{s'} \pmod{2}$, then equation 3.2 tells us that $\Psi_{s,s'}/2$ is the angle corresponding to one of the representations of $\prod_j p_j^{|\gamma_{j,s} - \gamma_{j,s'}|/2}$ as the sum of two squares since $\frac{t_s - t_{s'}}{8} = \frac{k}{4}$ for some $k \in \{0, 1, 2, 3\}$. Indeed, the only possibilities are $k = 0, 1$.
- If $t_s \not\equiv t_{s'} \pmod{2}$, then the correspondence is with a lattice point on $2 \prod_j p_j^{|\gamma_{j,s} - \gamma_{j,s'}|/2}$. This is because we can write $\frac{t_s - t_{s'}}{8} = \frac{1}{8} + \frac{k}{4}$ for $k = 0$ or $k = 1$.

This tells us that if $s \neq s'$, the quantity $2\pi\|\frac{\Psi_{s,s'}}{2}\|$ is the angle of a lattice point on the circle of radius

$$2^{\nu/2} \prod_j p_j^{|\gamma_{j,s} - \gamma_{j,s'}|/4},$$

where $\nu = 0$ or $\nu = 1$ depending on which case (from above) we're in. This observation allows us to construct our first bound on $\|\frac{\Psi_{s,s'}}{2}\|$. Since we know that this lattice point is not on the real axis, we know that the imaginary part must be at least 1 (or at most -1). This allows us to lower bound the length of the arc spanned by the angle $2\pi\|\frac{\Psi_{s,s'}}{2}\|$.

$$(3.3) \quad 2\pi \left\| \frac{\Psi_{s,s'}}{2} \right\| \cdot \sqrt{2} \prod_j p_j^{|\gamma_{j,s} - \gamma_{j,s'}|/4} > 1.$$

By our assumption that all $m + 1$ lattice points lie on an arc of length $\sqrt{2}R^\theta$, from the arc length formula we obtain an upper bound

$$(3.4) \quad 2\pi \left\| \frac{\Psi_{s,s'}}{2} \right\| \leq \frac{\sqrt{2}R^\theta}{2R} = \frac{R^{\theta-1}}{\sqrt{2}}.$$

Remark 3.5. The lower bound in Equation 3.3 is asymptotically sharp (as $R \rightarrow \infty$) and the upper bound in Equation 3.4 is sharp. These bounds are sharp in the sense that there exists collections of lattice points of the circle that attain equality (or attain equality in the limit $R \rightarrow \infty$) in the bound; however, these sets of lattice points are in general *not* the same sets. For example, for any $\theta > 0$, the sequence of lattice points $\{(n, \pm 1)\}$ produces a sequence of arc lengths that get arbitrarily close to the lower bound of Equation 3.3; however, these points will not be the lattice points that yield equality in Equation 3.4. This suggests that there may be some room for improvement in the proof of this theorem.

Proof. To see that the lower bound is sharp, consider the point $(n, 1)$. We know that this is a lattice point of the circle $x^2 + y^2 = n^2 + 1$. The length of the arc between the points $(n, 1)$ and $(\sqrt{n^2 + 1}, 0)$ is

$$\sqrt{n^2 + 1} \cdot \arcsin\left(\frac{1}{\sqrt{n^2 + 1}}\right)$$

and we can calculate the limit using L'Hopital's rule,

$$\lim_{n \rightarrow \infty} \frac{\arcsin\left(\frac{1}{\sqrt{n^2+1}}\right)}{\frac{1}{\sqrt{n^2+1}}} = \lim_{x \rightarrow 0} \frac{\arcsin(x)}{x} = \lim_{x \rightarrow 0} \frac{1}{\sqrt{1-x^2}} = 1.$$

Finally, it's important to note that there are indeed infinitely many products such that

$$2 \prod_{p_j \equiv 1(4)} p_j^{|\gamma_{j,s} - \gamma_{j,s'}|/2} = n^2 + 1.$$

To see this, we can show that every odd prime divisor of $n^2 + 1$ is congruent to 1 (mod 4). The proof proceeds as follows. Suppose that $p \mid n^2 + 1$ and that $p = 4k + 3$ for some integer k . This immediately implies that $p - 1 = 2(2k + 1)$ and $n^2 \equiv -1 \pmod{p}$. Since p is prime (\mathbb{Z}_p is a field), we get that $n^{p-1} \equiv 1 \pmod{p}$. Combining this with the two implications, we get

$$\begin{aligned} n^{p-1} \equiv 1 \pmod{p} &\iff (n^{2(2k+1)}) \equiv 1 \pmod{p} \\ &\iff (-1)^{2k+1} \equiv 1 \pmod{p} \\ &\iff -1 \equiv 1 \pmod{p} \end{aligned}$$

which is a contradiction since $p \neq 2$.

For the upper bound, it is sharp by our assumption that there are the points all lie on an arc of length $\sqrt{2}R^\theta$. \square

Continuing with the proof, we can multiply the upper and lower bounds over all possible pairs of s, s' (*i.e.*, $\frac{m(m+1)}{2}$ times); alternatively phrased, we can apply the upper and lower bounds to the quantity $\prod_{s,s'} \|\frac{\Psi_{s,s'}}{2}\|$ to yield the following equation.

$$(3.5) \quad R^{(\theta-1)(m(m+1)/2)} \geq \prod_{s,s'} \frac{1}{\prod_{p_j \equiv 1(4)} p_j^{|\gamma_{j,s} - \gamma_{j,s'}|/4}} = \prod_{p_j \equiv 1(4)} p_j^{-\sum_{s,s'} |\gamma_{j,s} - \gamma_{j,s'}|/4}.$$

Remark 3.6. The bound in Equation 3.5 is no longer sharp. It is only possible for one pair of lattice points s, s' to achieve the upper bound of Equation 3.4.

To make Equation 3.5 something more tractable, we seek to further lower bound the RHS by maximizing the sum $\sum_{s,s'} |\gamma_{j,s} - \gamma_{j,s'}|$.

However, before we continue, we address the choice of multiplying the bound to itself over all pairs of s, s' . This choice is somewhat arbitrary, but it yields something tractable. For example, if we consider the bound by itself, we get

$$R^{\theta-1} \geq \frac{1}{\prod_{p_j \equiv 1(4)} p_j^{|\gamma_{j,s} - \gamma_{j,s'}|/4}} \geq \frac{1}{\prod_{p_j \equiv 1(4)} p_j^{2\alpha_j/4}} = \frac{1}{R},$$

which doesn't produce anything helpful. We don't have to upper bound the difference $|\gamma_{s,j} - \gamma_{s',j}|$ though since the bound of Equation 3.4 is independent of s, s' . If we instead try to minimize the difference, we get

$$R^{\theta-1} \geq \frac{1}{\prod_{p_j \equiv 1(4)} p_j^{2/4}},$$

which also isn't very helpful. Ideally, we'd like a bound with an $R^{(\text{something})}$ on both sides so that we can say something explicit about θ . To get something of that form, we need to recover the α_j 's from the sum.

Returning to the task of maximizing $\sum_{s,s'} |\gamma_{j,s} - \gamma_{j,s'}|$, recall that $\gamma_{j,s} \in \{\gamma \in \mathbb{Z} \mid |\gamma| \leq \alpha_j, \gamma \equiv \alpha_j \pmod{2}\}$. This implies that the maximum value of a single term is $2\alpha_j$ and this happens when $\gamma_{j,s} = \alpha_j$ and $\gamma_{j,s'} = -\alpha_j$ (of course their labels can be switched).

Proposition 3.7.

$$\sum_{s,s'} |\gamma_{j,s} - \gamma_{j,s'}| \leq \alpha_j \frac{(m+1)^2 - \delta(m)}{2}$$

where $\delta(m) = 0$ if m is odd, and 1 if m is even.

Furthermore, this sum attains its maximum when $\lfloor (m+1)/2 \rfloor$ of the γ 's take a value of α_j and the remaining $\lceil (m+1)/2 \rceil$ take the value of $-\alpha_j$.

Proof. One can show that this distribution of values indeed achieves the maximum by considering what happens if we change the distribution in any way. For convenience, let $l = \lfloor (m+1)/2 \rfloor$.

Suppose that $\gamma_{j,1}, \dots, \gamma_{j,l} = \alpha_j$. If we change one of these, say $\gamma_{j,1}$ to any other value $\alpha_j - k =: \gamma'_{j,1}$, the terms of the sum that changes are those that involve $\gamma_{j,1}$ and these terms change as follows

$$\begin{aligned} \sum_{s'=2}^{m+1} |\gamma_{j,1} - \gamma_{j,s'}| - \sum_{s'=2}^{m+1} |\gamma'_{j,1} - \gamma_{j,s'}| &= 2l\alpha_j - \left(\sum_{s'=2}^l |\alpha_j - k - \alpha_j| + \sum_{s'=(l+1)}^{m+1} |\alpha_j - k + \alpha_j| \right) \\ &= 2l\alpha_j - (l-1)k - l(2\alpha_j - k) \\ &= k \end{aligned}$$

which shows us that the altered sum is less than the original. Another way of viewing this is to take the perspective of devising a greedy algorithm to pick values of γ_j with the goal of maximizing the sum in the statement of the proposition. We have proven that this greedy algorithm yields the optimal selections by showing that every step it takes is the best possible.

Now, we need to calculate the value of the sum given that $\lfloor (m+1)/2 \rfloor$ of the γ 's are α_j and the remaining $\lceil (m+1)/2 \rceil$ are $-\alpha_j$. Like before, we can assume that the first $\lfloor (m+1)/2 \rfloor$ are α_j .

$$\sum_{s,s'} |\gamma_{j,s} - \gamma_{j,s'}| = \sum_{s=1}^m \sum_{s'=s+1}^{m+1} |\gamma_{j,s} - \gamma_{j,s'}|$$

$$\begin{aligned}
&= \sum_{s=1}^{\lfloor (m+1)/2 \rfloor} \sum_{s'=s+1}^{m+1} |\alpha_j - \gamma_{j,s'}| + \sum_{s=\lfloor (m+1)/2 \rfloor + 1}^{m+1} \sum_{s'=s+1}^{m+1} |-\alpha_j - \gamma_{j,s'}| \\
&= \left\lfloor \frac{m+1}{2} \right\rfloor \sum_{s'=\lfloor (m+1)/2 \rfloor + 1}^{m+1} |\alpha_j - (-\alpha_j)| + \left\lfloor \frac{m+1}{2} \right\rfloor \sum_{s'=\lfloor (m+1)/2 \rfloor + 1}^{m+1} |-\alpha_j - (-\alpha_j)| \\
&= \left\lfloor \frac{m+1}{2} \right\rfloor \left\lfloor \frac{m+1}{2} \right\rfloor (2\alpha_j) + 0 \\
&= \begin{cases} \frac{(m+1)^2}{2} \alpha_j & m \equiv 1 \pmod{2} \\ \frac{(m+1)^2 - 1}{2} \alpha_j & m \equiv 0 \pmod{2} \end{cases}
\end{aligned}$$

□

Remark 3.8. Proposition 3.7 is asymptotically sharp (as $R \rightarrow \infty$). Like before, we are using the term sharp in the sense that there for any m , there exists a circle with lattice points that satisfy the condition of the Proposition (*i.e.*, yield equality in the upper bound of the sum). Importantly, these points may not lie on the arc of length $\sqrt{2}R^\theta$ and the set of points that satisfies equality for one p_j do not in general satisfy equality for other p_j . One can expect that this bound can be improved (decreased) by taking into consideration the requirement that all $m+1$ lattice points lie on this arc of interest.

Proof. Recall that we are working with a circle centered at the origin of radius $R = \sqrt{n}$ where

$$n = \prod_{p_j \equiv 1(4)} p_j^{\alpha_j}$$

and that the lattice points on this circle are parameterized by the formula (Equation 3.2)

$$\sqrt{n} e^{2\pi i (\sum_j \gamma_j \Phi_j + t/4)}$$

where $\gamma_j \in \Lambda_{\alpha_j}$ and $t \in \{0, 1, 2, 3\}$, and that this parameterization shows that there are $4 \prod_j (1 + \alpha_j)$ lattice points on the circle. Furthermore, it is clear that each lattice point λ_s corresponds to unique set of parameters $\{\gamma_{j,s}\}_{p_j \equiv 1(4)} \cup \{t_s\}$, and since there are the same amount of lattice points as parameter sets, we do have a 1-1 correspondence between lattice points and the parameters that define them. We ask the question: For each p_j , is it possible for there to be $\lfloor \frac{m+1}{2} \rfloor$ lattice points that with $\gamma_j = \alpha_j$?

For fixed m , the answer depends on the radius via the number of prime factors p_j .

For example purposes, fix p_j for now. If $m = 1$, then there is certainly one lattice point with $\gamma_j = \alpha_j$. Similarly, if $m = 3$, there will be two lattice points with $\gamma_j = \alpha_j$ separated by an angle of $t\pi/2$ where $t = 1, 2, 3$. However, if we consider a radius with only one prime factor p_j and $m = 10$, then there can not be 5 lattice points with $\gamma_j = \alpha_j$.

Since we can count the number of lattice points with $\gamma_j = \alpha_j$ by fixing the value and letting the other parameters range over their respective ranges $\Lambda_{\alpha'_j}$, we see that in general, there will be at most

$$4 \prod_{\substack{p_i \equiv 1(4) \\ p_i \neq p_j}} (1 + \alpha_i)$$

lattice points with $\gamma_j = \alpha_j$ and if the product is empty we say it takes the value 1. Furthermore, the above count also applies if α_j is replaced with any other element of Λ_{α_j} . However, for any $m+1$, since the number of lattice points on a circle centered at the origin is unbounded as $R \rightarrow \infty$, we can always find R sufficiently large such that it is possible for there to be $\lfloor (m+1)/2 \rfloor$ lattice points with $\gamma_j = \alpha_j$. □

We will look at another example to make this remark even clearer.

Example 3.9. We again consider the circle $x^2 + y^2 = 5^2$. In this case $\alpha = 2$, so $\Lambda_\alpha = \{-2, 0, 2\}$. Recall from earlier that 5 splits into $(2 + i)(2 - i)$ in $\mathbb{Z}[i]$ and by our convention of taking the point in the first octant, we have $\omega = 2 + i$. Then, the lattice points corresponding to the choices that maximize the sum in the Proposition are the points: $\omega^2, \bar{\omega}^2, i\omega^2, i\bar{\omega}^2, -\omega^2, -\bar{\omega}^2, -i\omega^2, -i\bar{\omega}^2$. Geometrically, these points are the 8 lattice points on the circle that do not lie on the x or y axis (*i.e.*, the points $(\pm 3, \pm 4)$ and $(\pm 4, \pm 3)$). This tells us that indeed the lattice points that maximize the sum can be rather far apart. For example if $m + 1 = 3$, Proposition 3.7 tells us that we want 2 points that correspond to $\alpha = 2$ and one that corresponds to $\alpha = -2$. However, just by referencing the list above, we see that these three lattice points span an arc that is *at least* one quarter of the circle.

For another example, consider the circle $x^2 + y^2 = 5 \cdot 13$. In this case, $\alpha_5 = \alpha_{13} = 1$, so the sets $\Lambda_{\alpha_5} = \Lambda_{\alpha_{13}} = \{-1, 1\}$. Since there are only two possible values for α_j , all 16 of the lattice points on this circle are eligible candidates for members of the set of $m + 1$ lattice points that maximizes the sum in Proposition 3.7. This example shows that there are also circles where the lattice points that satisfy the maximum of Proposition 3.7 can also be somewhat close together. In general, if $n = \prod_{p_j \equiv 1 \pmod{4}} p_j$ and we have a circle of radius $r = \sqrt{n}$, then all the lattice points on this circle will be eligible candidates for the maximum-achieving selection.

Lastly, consider the circle $x^2 + y^2 = 5^2 \cdot 13$. In this case, $\Lambda_{\alpha_5} = \{-2, 0, 2\}$ and $\Lambda_{\alpha_{13}} = \{-1, 1\}$. In Figure 2 we see that if $m + 1 = 4$ then we can pick lattice points that are as close together as possible; however, if $m + 1 > 4$, then the fifth lattice point we choose for the set that achieves the maximum of Proposition 3.7 is no longer the closest possible choice. Notice that the lattice points $(\pm 10, \pm 15)$ and $(\pm 15, \pm 10)$ are missing from Figure 2. Just to drive the point home, if we carry out the same exercise for $n = 5^3 \cdot 65$, the points with $\alpha_5 = \pm 3$ look like every-other lattice point on the circle.

Returning to the proof of the Theorem, by applying Proposition 3.7 to our current bound (Equation 3.5), we get

$$R^{(\theta-1)(m(m+1)/2)} > \left(\prod_{p_j \equiv 1(4)} p_j^{\alpha_j \frac{(m+1)^2 - \delta(m)}{4}} \right)^{-1/2} = R^{-\frac{(m+1)^2 - \delta(m)}{4}}$$

and we calculate that

$$\begin{aligned} (\theta - 1) \left(\frac{m(m+1)}{2} \right) &> -\frac{(m+1)^2 - \delta(m)}{4} \\ \theta \left(\frac{m(m+1)}{2} \right) &> -\frac{(m+1)^2 - \delta(m)}{4} + \frac{m(m+1)}{2} \\ \theta &> 1 - \frac{(m+1)^2 - \delta(m)}{2m(m+1)} = \begin{cases} 1 - \frac{(m+1)}{2m} & m \text{ odd} \\ 1 - \frac{m+2}{2(m+1)} & m \text{ even} \end{cases} \\ &= 1 - \frac{(m+1) + \delta(m)}{2(m + \delta(m))} \\ &= \frac{m - 1 + \delta(m)}{2(m + \delta(m))} \\ &= \frac{1}{2} - \frac{1}{2m + 2\delta(m) \pm 2} \end{aligned}$$

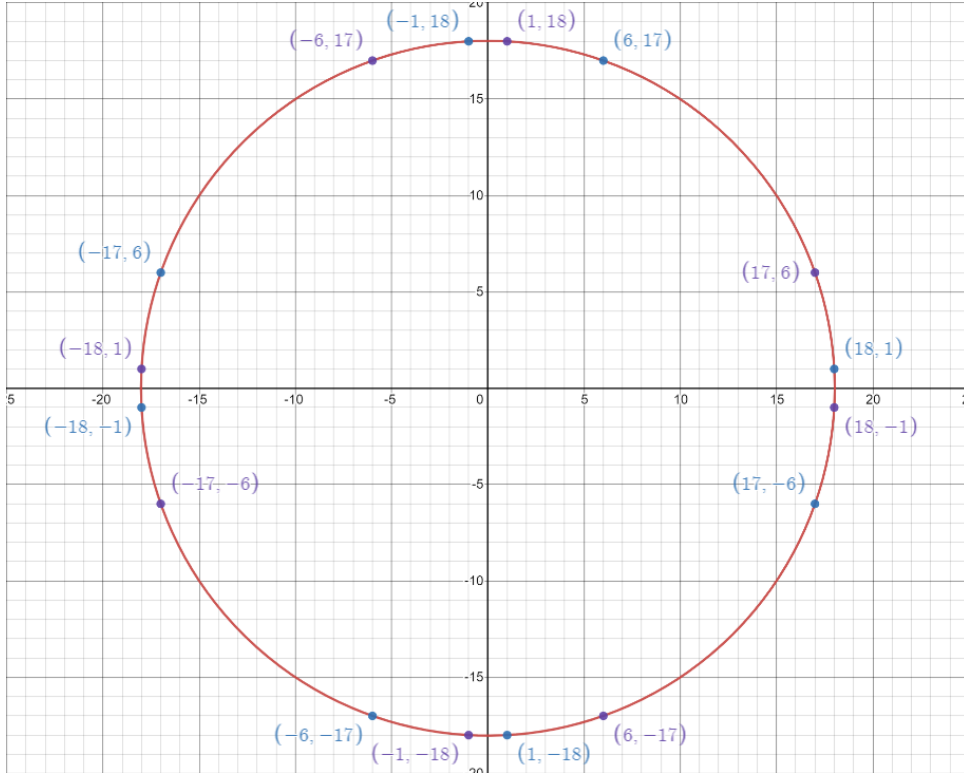


FIGURE 2. Lattice points on the circle $x^2 + y^2 = 5^2 \cdot 13$ with $\alpha_5 = \pm 2$. The purple lattice points correspond to $\alpha_5 = 2$ and the blue to $\alpha_5 = -2$.

$$= \frac{1}{2} - \frac{1}{4 \lfloor m/2 \rfloor + 2}$$

completing the proof of Theorem 3.1. We have also shown that most steps of this proof is optimal. By Remarks 3.5 and 3.8 we know that each bound used in the proof is sharp. However, as mentioned in Remark 3.6, when multiplying the bound of $\|\Psi_{s,s'}/2\|$ over all pairs of s, s' , the bound is no longer sharp since only one pair of s, s' will achieve equality of the upper bound. In Section 4 we will discuss some ways that the result may be improved.

3.3. Open questions.

So far, we have a bound for arcs of the form R^θ where $\theta \in [0, 1/2)$. What happens when $\theta \in [1/2, 1)$? We don't know! What about the case where $\theta = 1$? Thankfully, for this case, we know there can not be a uniform bound. The intuition is that any arc with length R will always contain a non-trivial slice of the circle as $R \rightarrow \infty$. On the contrary, an arc of length $R^{1/2}$ represents a smaller and smaller piece of the pie as $R \rightarrow \infty$ (more generally, we can replace $1/2$ with $1 - \varepsilon$ for any $\varepsilon > 0$).

Proposition 3.10. *For any arc of length $\ll R$, the maximum number of lattice points on such an arc is unbounded in R .*

Proof. For any R , we can write the length of the arc as CR for some constant C . Suppose that on any arc of length CR , there exists a finite upper bound M such that there are at most M lattice points on the arc. However, since $\exists k$ that does *not* depend on R such that $kCR \geq 2\pi R$, we can upper bound the number of lattice points on the circle by kM . In other words, this implies that $r(R^2) \leq kM$ for all R where $r(n)$ is the number of lattice

points on the circle $x^2 + y^2 = n$. However, this is a contradiction since we know $r(R^2)$ is unbounded in R . \square

For the case of $\theta \in [1/2, 1)$, the following conjecture by Cilleruelo and Granville [CG09] states:

Conjecture 3.11. For every $\varepsilon > 0$, there exists a constant $M_\varepsilon < \infty$ such that there are no more than M_ε lattice points on an arc of length $R^{1-\varepsilon}$ on a circle of radius R that is centered at the origin.

Furthermore, under this conjecture, M_ε must be at least $e^{C/\varepsilon}$ for some constant $C > 0$.

Alongside the conjecture, Cilleruelo and Granville provide the following construction (details added here are not included in their paper).

Fix some positive integer m . Define the choice function $\sigma : \{1, \dots, 2m\} \rightarrow \{-1, 1\}$. In other words, for every integer up to $2m$, we assign it either 1 or -1 . Notice that there are then exactly $\binom{2m}{m}$ such σ where $\sum_{j=1}^{2m} \sigma(j) = 0$. Now, we define

$$\Sigma_m := \left\{ \sigma : \{1, \dots, 2m\} \rightarrow \{-1, 1\} \mid \sum_{j=1}^{2m} \sigma(j) = 0 \right\}$$

and let $l = 1, \dots, \binom{2m}{m}$ index this set.

Let a be some very large integer; it is difficult to say precisely how large a must be, but it must be at least large enough such that $a + 2m \approx a$. We now define the gaussian integer

$$v_l = \prod_{j=1}^{2m} (a + j + i\sigma_l(j))$$

for $l = 1, \dots, \binom{2m}{m}$. Now, notice that all of these gaussian integers have the same norm because they are all the product of $2m$ gaussian integers of the form $a + j \pm i$. This tells us they all lie on the circle of radius $R = \prod_{j=1}^{2m} |a + j + i| = \mathcal{O}(a^{2m})$. Furthermore, it's also important to note that each component of the product in the definition of v_l has imaginary part equal to ± 1 . This tells us that v_l is a product of many Gaussian integers with very small argument, and since the argument of the product is the sum of the arguments, intuitively, we should expect that the v_l will not end up too far from the real line. In fact, we can calculate that

$$\begin{aligned} v_l - v_k &= \prod_{j=1}^{2m} (a + j + i\sigma_l(j)) - \prod_{j=1}^{2m} (a + j + i\sigma_k(j)) \\ &= \left(\prod_{j=1}^{2m} (a + j) + a^{2m-1} \sum_{j=1}^{2m} i\sigma_l(j) \right. \\ &\quad \left. + a^{2m-2} \sum_{1 \leq j < j' \leq 2m} (j i \sigma_l(j') + j' i \sigma_l(j) + i^2 \sigma_l(j) \sigma_l(j')) + \mathcal{O}(a^{2m-3}) \right) \\ &\quad - \left(\prod_{j=1}^{2m} (a + j) + a^{2m-1} \sum_{j=1}^{2m} i\sigma_k(j) \right. \\ &\quad \left. + a^{2m-2} \sum_{1 \leq j < j' \leq 2m} (j i \sigma_k(j') + j' i \sigma_k(j) + i^2 \sigma_k(j) \sigma_k(j')) + \mathcal{O}(a^{2m-3}) \right) \end{aligned}$$

$$\begin{aligned}
&= a^{2m-2} \sum_{1 \leq j < j' \leq 2m} (j(i\sigma_l(j') - i\sigma_k(j')) + j'(i\sigma_l(j) - i\sigma_k(j))) \\
&\quad + i^2(\sigma_l(j)\sigma_l(j') - \sigma_k(j)\sigma_k(j')) + \mathcal{O}(a^{2m-3})
\end{aligned}$$

Thus, to bound the difference between the lattice points, we want to bound this double sum. By similar reasoning to the proof of Proposition 3.7, this sum is maximized when

$$\sigma_l(j) - \sigma_k(j) = \begin{cases} -2 & \text{for } j = m+1, m+2, \dots, 2m \\ 2 & \text{for } j = 1, 2, \dots, m \end{cases}.$$

Working in isolation from the rest of the terms, we can rewrite/manipulate the first two components of the sum as follows.

$$\begin{aligned}
\cdots &= \sum_{j=1}^{2m-1} \sum_{j'=j+1}^{2m} \left(j'(i\sigma_l(j) - i\sigma_k(j)) + j(i\sigma_l(j') - i\sigma_k(j')) \right) \\
&= \sum_{j=1}^m \sum_{j'=j+1}^{2m} j'(2i) + \sum_{j=m+1}^{2m-1} \sum_{j'=j+1}^{2m} j'(-2i) \\
&\quad + \sum_{j=1}^{m-1} \left(\sum_{j'=j+1}^m j(2i) + \sum_{j'=m+1}^{2m} j(-2i) \right) + \sum_{j=m}^{2m-1} \sum_{j'=j+1}^{2m} j(-2i) \\
&= 2i \left(\sum_{j=1}^m \sum_{j'=j+1}^{2m} j' - \sum_{j=m+1}^{2m-1} \sum_{j'=j+1}^{2m} j' \right) \\
&\quad + \sum_{j=1}^{m-1} j(2i(m-j) - 2im) - \sum_{j=m}^{2m-1} j(2i(2m-j)) \\
&= i \left(\sum_{j=1}^m (2m-j)(2m+j+1) - \sum_{j=m+1}^{2m-1} (2m-j)(2m+j+1) \right) \\
&\quad - 2i \sum_{j=1}^{m-1} j^2 + 2i \sum_{j=m}^{2m-1} j^2 - 4im \sum_{j=m}^{2m-1} j \\
&= i(2m^3 + 2m^2) + 2i(-m^2 + 2m^3) + 4i(-3m^3/2 + m^2/2) \\
&= 2im^2
\end{aligned}$$

which tells us that the double sum in the original expression is $\mathcal{O}(m^2)$. In other words, we have that

$$\begin{aligned}
|v_l - v_k| &\ll a^{2m-2} m^2 \\
&\ll_m R^{1-1/m}
\end{aligned}$$

where the implicit constant depends on m in the way of m^2 . This bound implies that all $\binom{2m}{m}$ lattice points lie on an arc of length $\mathcal{O}(R^{1-1/m})$. To explain the final claim of the conjecture, we can use Stirling's approximation to calculate that

$$\binom{2m}{m} = \frac{(2m)!}{m!m!} \sim \frac{\sqrt{2\pi(2m)} \left(\frac{2m}{e}\right)^{2m}}{\sqrt{2\pi m} \left(\frac{m}{e}\right)^m \sqrt{2\pi m} \left(\frac{m}{e}\right)^m} = \frac{2^{2m}}{\sqrt{\pi m}} = \mathcal{O}(2^{2m}).$$

Intuitively, this makes sense as choosing m elements from $2m$ is a special case of choosing n elements from a total of $2m$, and this is equivalent to assigning a “yes” or “no” to each of the $2m$ elements. This tells us that the uniform bound M_ε on the number of lattice points on an arc of length $\mathcal{O}(R^{1-\varepsilon})$ must be $\mathcal{O}(2^{1/\varepsilon})$ which can be written in the form $e^{C/\varepsilon}$ for some constant $C > 0$.

Example 3.12. As an example, we can compute what happens in the case of $m = 1$. In this case, the construction above promises us 2 lattice points on an arc of length $\mathcal{O}(1)$. Since $m = 1$, the only two choice functions σ are

$$\sigma_1 : \{1, 2\} \rightarrow \{1, -1\} \quad \text{and} \quad \sigma_2 : \{1, 2\} \rightarrow \{-1, 1\}$$

and correspondingly,

$$\begin{aligned} v_1 &= (a + 1 + i)(a + 2 - i) & \text{and} & \quad v_2 = (a + 1 - i)(a + 2 + i) \\ &= a^2 + 3a + 3 + i & & \quad = a^2 + 3a + 3 - i \end{aligned}$$

which are always two lattice points that are always $< \pi$ units apart.

4. A POTENTIAL IMPROVEMENT ON THEOREM 3.1

In the discussion of the proof of Theorem 3.1, we pointed out a few places where the bounds used by Cilleruelo and Cordoba are not sharp. Naturally, we wonder if it possible to improve their result by improving these bounds. If so, we might be able to say something about the case where $\theta \in [1/2, 1)$.

For legibility, the important bounds from the previous section are repeated below:

$$\frac{1}{\sqrt{2}} \prod_j p_j^{-|\gamma_{j,s} - \gamma_{j,s'}|/4} < 2\pi \left\| \frac{\Psi_{s,s'}}{2} \right\| \leq \frac{R^{\theta-1}}{\sqrt{2}}.$$

Repeating what was said before, the main loss in their bounds comes from the fact that they multiply the above sequence of inequalities over all $m(m+1)/2$ pairs of lattice points s, s' . The issue is that only 1 of the pairs can attain equality of the above upper bound, and for m somewhat large we can intuitively expect that there will be some lattice points that are much closer together than the bound above. It is difficult to be precise about what “somewhat large” means, but if we consider the example of $m = 1, 2$, then we see that this upper bound is not wrong by that much. If $m = 1$, then $m + 1 = 2$ and there is only one possible pair of lattice points, so there is no loss at all. If $m = 2$, then $m + 1 = 3$ and if the middle lattice point is not too close to either of the sides, then the product is still somewhat accurate although it will be off by a constant factor, and this can be seen in a comparison of Theorem 3.1 and Theorem 2.2.

In trying to improve this result, if we assume that $\theta \leq 1$, then we arrive at an upper bound for $\|\Psi_{s,s'}/2\|$ that depends on s, s' . First, define

$$r_{s,s'} := \prod_j p_j^{\frac{|\gamma_{j,s} - \gamma_{j,s'}|}{4}}.$$

Recall that $R = \prod_j p_j^{2\alpha_j/4}$. Since $\frac{r_{s,s'}}{R} \leq 1$ for all s, s' and since $\theta \leq 1$, it follows that

$$\frac{r_{s,s'}}{R} \leq \left(\frac{r_{s,s'}}{R} \right)^\theta \quad \implies \quad 1 \leq \left(\frac{r_{s,s'}}{R} \right)^{\theta-1}.$$

Then, we can multiply the original upper bound on the LHS by the LHS of the above and on the RHS by the RHS of the above to get

$$2\pi \left\| \frac{\Psi_{s,s'}}{2} \right\| \leq \frac{r_{s,s'}^{(\theta-1)}}{\sqrt{2}} = \frac{1}{\sqrt{2}} \prod_j p_j^{\frac{|\gamma_{j,s} - \gamma_{j,s'}|}{4}(\theta-1)}.$$

If we then multiply over all pairs s, s' , we get

$$\prod_{s,s'} 2\pi \sqrt{2} \left\| \frac{\Psi_{s,s'}}{2} \right\| \leq \prod_j p_j^{(\theta-1) \sum_{s,s'} \frac{|\gamma_{j,s} - \gamma_{j,s'}|}{4}}.$$

Since $\theta - 1 \leq 0$, to maximize this product, we need to minimize the sum in the exponent. We can compute that the average value of $|\gamma - \gamma'|$ is

$$\begin{aligned} \frac{1}{\alpha(\alpha+1)/2} \sum_{\gamma, \gamma' \in \Lambda_\alpha} |\gamma - \gamma'| &= \frac{1}{\alpha(\alpha+1)/2} \sum_{k=0}^{\alpha-1} \sum_{l=k+1}^{\alpha} (2l - 2k) \\ &= \frac{2(\alpha+2)}{3}. \end{aligned}$$

The approach taken for this computation is that before computing all of the pairwise differences, we can first assume that the array is sorted in increasing order (*i.e.*, $\gamma \geq \gamma'$). Furthermore, the sum of the pairwise differences is invariant under shifting, so we can shift the whole array to be positive numbers to make it easier to work with; this allows us to write all the elements in the form $2\alpha - 2k$ for $k = 0 \dots \alpha$.

It's unclear if computing the average is helpful, but perhaps it will be a useful tool in trying to bound the sum of differences.

This new upper bound also shows us that the lower bound used in Cilleruelo and Cordoba's argument corresponds to the case $\theta = 0$. To recap, if we define

$$U_m := \max_{s,s'} \sum |\gamma_s - \gamma_{s'}| \quad \text{and} \quad L_m := \min_{s,s'} \sum |\gamma_s - \gamma_{s'}|,$$

where the maximum and minimum are taken over the set of all possible collections of $m+1$ lattice points that lie on an arc of length $\sqrt{2}R^\theta$. Inserting this into our current bounds yields

$$\prod_j p_j^{-U_m} \leq \prod_j p_j^{-\sum_{s,s'} \frac{|\gamma_{j,s} - \gamma_{j,s'}|}{4}} < \prod_j p_j^{(\theta-1) \sum_{s,s'} \frac{|\gamma_{j,s} - \gamma_{j,s'}|}{4}} \leq \prod_j p_j^{(\theta-1)L_m}$$

which in turn would imply

$$\theta > 1 - \frac{U_m}{L_m}.$$

We would like expressions for U_m and L_m that are of the same order (or for the order of L_m to be larger) with respect to m so that the bound remains sensible as $m \rightarrow \infty$. However, this approach may not actually yield anything useful. One problem is that we would like $U_m/L_m \in (0, 1)$ for all m for this bound to make sense, but this seems problematic because why should the ratio of the upper bound to the lower bound of the same quantity lie in the interval $(0, 1)$? It's also reasonable to desire that as $m \rightarrow \infty$, the quantity $1 - U_m/L_m \rightarrow 1$.

Lastly, recall that in Cilleruelo and Cordoba's proof they used the bound

$$\sum_{s,s'} |\gamma_{j,s} - \gamma_{j,s'}| \leq \alpha_j \frac{(m+1)^2 - \delta(m)}{2}$$

but that the lattice points that satisfy equality of this bound can actually be rather far apart (Example 3.9). In order to improve the Theorem by improving steps of the proof, we would likely have to improve this upper bound since it seems very difficult to find a lower bound of the sum that will yield a meaningful result.

5. LATTICE POINTS WITHIN A NEIGHBORHOOD

In this section, we consider the problem of counting the number of lattice points within a neighborhood of a curve. Although we are generally interested in the 2-dimensional case, the content of subsection 5.2 lives in general k -dimensional space.

5.1. A close neighborhood of an ellipse.

We now turn our attention to a Lemma by Bourgain [Bou09] in which he bounds the number of lattice points within a small neighborhood of a “regular” ellipse. The proof of this lemma follows a similar argument to that of Theorem 2.2.

Lemma 5.1 ([Bou09]).

Let \mathcal{E} be a regular ellipse of size R . Denote by \mathcal{E}_ε the ε -neighborhood of \mathcal{E} . Then,

$$\max_{a \in \mathbb{R}^2} \#\{B(a, R^{1/3}) \cap \mathcal{E}_{1/R} \cap \mathbb{Z}^2\} < C$$

and, in particular,

$$\#\{\mathcal{E}_{1/R} \cap \mathbb{Z}^2\} < CR^{2/3}.$$

In other words, there are at most C lattice points within a small neighborhood of a section of the ellipse that has size $R^{1/3}$, and the total number of integer lattice points in $\mathcal{E}_{1/R}$ is at most $CR^{2/3}$ for some constant C .

In his paper, Bourgain does not define what a “regular ellipse of size R ” is. We will define an ellipse of size R as an ellipse of the form

$$\frac{x^2}{(kR)^2} + \frac{y^2}{((1-k)R)^2} = 1$$

for $k \in (0, 1)$. which is an ellipse centered at the origin with major/minor axes parallel to the coordinate axes. We choose this definition since it is consistent with the concept of a circle of “size” (diameter) R . Another perspective on this definition is that ellipses of the form described above have constant eccentricity as $R \rightarrow \infty$. Eccentricity is a way of measuring how much a given conic section deviates from the circle which has eccentricity 0. The formula for eccentricity of an ellipse with major/minor axes lengths a, b ($a > b$) is

$$e = \sqrt{1 - \frac{b^2}{a^2}}.$$

Notice that if a and b do not grow at the same rate in R , then as $R \rightarrow \infty$, the eccentricity will converge to 1 and the ellipse will become a parabola. As we will see later, this definition is also essential to the truth of Bourgain’s lemma (Remark 5.2).

Bourgain’s proof starts by supposing there exists P_1, P_2 , and P_3 three noncolinear lattice points in $B(a, cR^{1/3}) \cap \mathcal{E}_{1/R}$ where c is some sufficiently small number to be determined later. Since these points are not on the same line, the triangle they create has non-zero area. Bourgain then writes,

$$0 \neq \text{area triangle}(P_1, P_2, P_3) = \frac{1}{2} \left| \begin{array}{ccc} 1 & 1 & 1 \\ P_1 & P_2 & P_3 \end{array} \right| \in \frac{1}{2}\mathbb{Z}_+.$$

Bourgain has provided some rather cryptic notation for the area of a triangle! It appears that the expression within and including the bars is meant to represent the determinant of a 3×3 matrix with P_1 denoting the column vector $(P_{1,x}, P_{1,y})^\top$ where P_x, P_y are the x, y coordinates of the lattice point. This is because the area of a triangle with vertices A, B, C can be written

$$\begin{aligned} \text{area triangle}(A, B, C) &= \frac{1}{2} |A_x(B_y - C_y) + B_x(C_y - A_y) + C_x(A_y - B_y)| \\ &= \frac{1}{2} |(A_x B_y - A_y B_x) + (B_x C_y - B_y C_x) - (A_x C_y - A_y C_x)| \\ &= \frac{1}{2} \left| \det \begin{bmatrix} 1 & 1 & 1 \\ A_x & B_x & C_x \\ A_y & B_y & C_y \end{bmatrix} \right| \end{aligned}$$

Notice that there is an absolute value in this expression but there is not in Bourgain's equation; this is ok since we can choose the labels of the points so that the determinant is positive. Another way to interpret this expression is to view the matrix as a change of basis. The determinant of this matrix tells us how the area of the unit square gets scaled, and then half of that gives us the area of the triangle.

Continuing, like before, since P_1, P_2, P_3 are lattice points, we know that

$$\text{area triangle}(P_1, P_2, P_3) \geq 1/2.$$

Bourgain continues by defining three more points P'_1, P'_2, P'_3 that lie *on* the ellipse \mathcal{E} . Importantly, these points do not necessarily lie on the integer lattice. We choose the points such that $\|P_j - P'_j\| < 1/R$ for $j = 1, 2, 3$. Here, and later, $\|x - y\|$ denotes the Euclidean distance between the points x, y . Next, Bourgain wants to bound the difference (written in his notation)

$$\left| \left| \begin{array}{ccc} 1 & 1 & 1 \\ P_1 & P_2 & P_3 \end{array} \right| - \left| \begin{array}{ccc} 1 & 1 & 1 \\ P'_1 & P'_2 & P'_3 \end{array} \right| \right|.$$

A bit of calculation shows that

$$\begin{aligned} &= \left| (P_{1,x}P_{2,y} - P_{2,x}P_{1,y}) + (P_{2,x}P_{3,y} - P_{3,x}P_{2,y}) + (P_{3,x}P_{1,y} - P_{1,x}P_{3,y}) \right. \\ &\quad \left. - ((P'_{1,x}P'_{2,y} - P'_{2,x}P'_{1,y}) + (P'_{2,x}P'_{3,y} - P'_{3,x}P'_{2,y}) + (P'_{3,x}P'_{1,y} - P'_{1,x}P'_{3,y})) \right| \\ &= \left| [(P_{1,x}P_{2,y} - P_{2,x}P_{1,y}) - (P'_{1,x}P'_{2,y} - P'_{2,x}P'_{1,y})] \right. \\ &\quad \left. + [(P_{2,x}P_{3,y} - P_{3,x}P_{2,y}) - (P'_{2,x}P'_{3,y} - P'_{3,x}P'_{2,y})] \right. \\ &\quad \left. + [(P_{3,x}P_{1,y} - P_{1,x}P_{3,y}) - (P'_{3,x}P'_{1,y} - P'_{1,x}P'_{3,y})] \right| \\ &= \left| (P_{1,x} - P_{3,x})(P_{2,y} - P'_{2,y}) + (P_{3,x} - P_{2,x})(P_{1,y} - P'_{1,y}) + (P_{2,x} - P_{1,x})(P_{3,y} - P'_{3,y}) \right. \\ &\quad \left. + (P'_{3,y} - P'_{1,y})(P_{2,x} - P'_{2,x}) + (P'_{1,y} - P'_{2,y})(P_{1,x} - P'_{1,x}) + (P'_{2,y} - P'_{3,y})(P_{3,x} - P'_{3,x}) \right| \\ &\leq 3(2cR^{1/3} \frac{1}{R}) + 3((2cR^{1/3} + 2R^{-1}) \frac{1}{R}) < 6cR^{1/3}R^{-1} + \frac{1}{2}R^{1/3}R^{-1} < R^{1/3}R^{-1} \end{aligned}$$

Earlier we asked that c be “sufficiently small” and here we see that we require $c < 1/12$ (importantly, c does not depend on R). The rest of the inequalities come from the facts that the points $P_1, P_2, P_3 \in B(a, cR^{1/3})$ and since $\|P'_{j,x} - P'_{k,x}\| \leq 2cR^{1/3} + 2/R < \frac{1}{2}R^{1/3}$ for any $R > 4$. This bound implies that as $R \rightarrow \infty$, the difference between the areas of the triangles $P_1P_2P_3$ and $P'_1P'_2P'_3$ approaches zero. Furthermore, for $R > 4$, we see that

$$\frac{1}{2} \left| \begin{vmatrix} 1 & 1 & 1 \\ P_1 & P_2 & P_3 \end{vmatrix} - \begin{vmatrix} 1 & 1 & 1 \\ P'_1 & P'_2 & P'_3 \end{vmatrix} \right| < \frac{1}{2} 4^{-2/3} < \frac{1}{4}$$

which, in conjunction with the fact that $\text{area triangle}(P_1, P_2, P_3) \geq 1/2$, implies that

$$\text{area triangle}(P'_1, P'_2, P'_3) > 1/4.$$

This is equation (1.4.6) of Bourgain’s paper.

At the same time, we can upper bound the area of the triangle $P'_1P'_2P'_3$ by using the formula from that states that the area = $abc/4r$ where r is the circumradius of the triangle with side-lengths a, b, c . The first thing to notice is that the ellipse centered at the origin, with major/minor axes parallel to the integer lattice, and of size R is contained in the circle centered at the origin of radius $R/2$; furthermore, we can transform the ellipse to the circle by scaling along the axis parallel to the minor axis. Thus, if Δ is a triangle on the ellipse \mathcal{E} , then the area of Δ is upper bounded by the area of the transformed triangle on the circle of radius $R/2$. This yields the following upper bound. $\|x - y\|$ denotes the Euclidean distance between the points $x, y \in \mathbb{R}^2$.

$$\begin{aligned} \text{area triangle}(P'_1, P'_2, P'_3) &\leq \frac{\|P'_1 - P'_2\| \|P'_2 - P'_3\| \|P'_1 - P'_3\|}{4(R/2)} \\ &\leq \frac{(2cR^{1/3} + 2R^{-1})^3}{2R} \\ &< cR^{1/3} \frac{R^{2/3}}{R} = c \end{aligned}$$

by using a similar bound to the previous calculation. However, this yields a contradiction since $c < \frac{1}{4}$ so the area can not simultaneously satisfy both inequalities. This shows there are either at most 2 integer lattice points in $B(a, cR^{1/3}) \cap \mathcal{E}_{1/R}$ or that the three points P_1, P_2, P_3 are co-linear.

If there are at most 2 integer lattice points in $B(a, cR^{1/3}) \cap \mathcal{E}_{1/R}$, then we know that $\exists C_1 > 0$ such that $\#\{B(a, R^{1/3}) \cap \mathcal{E}_{1/R}\} < C_1$ since we cover $B(a, R^{1/3})$ with finitely many balls of radius $cR^{1/3}$ and since a is the point that maximizes the number of lattice points in the intersection of the ball and the neighborhood of the ellipse.

If the points are co-linear, let Λ denote the line through them. For any ellipse of size R , we can write the equation of the ellipse in the form

$$\frac{x^2}{(kR)^2} + \frac{y^2}{((1-k)R)^2} = 1$$

for some $k \in (0, 1)$. We can also assume that the major axis is always in the x direction which means we assume WLOG that $k \geq 1/2$. In the case that $k = 1/2$, the ellipse is a circle of radius $R/2$. Then, the formulas for the outer and inner boundaries of the neighborhood can be written as

$$\frac{x^2}{(kR + R^{-1})^2} + \frac{y^2}{((1-k)R + R^{-1})^2} = 1 \quad \text{and} \quad \frac{x^2}{(kR - R^{-1})^2} + \frac{y^2}{((1-k)R - R^{-1})^2} = 1.$$

We want to show that for any Λ , the length of the intersection $\Lambda \cap \mathcal{E}_{1/R}$ is uniformly bounded in R . If this is the case, then we know that there are at most finitely many

lattice points in $\Lambda \cap \mathcal{E}_{1/R}$. For any fixed k , the longest possible line through $\mathcal{E}_{1/R}$ is the line parallel to the major axis that is tangent to the inner boundary. This line can be written as $y = (1 - k)R - R^{-1}$ and we can calculate that the x -coordinates of the intersection points with the outer boundary are

$$x = \pm 2 \sqrt{\frac{(1 - k)(kR^2 + 1)^2}{((1 - k)R^2 + 1)^2}}$$

by solving the quadratic. This tells us that

$$\text{len}(\Lambda \cap \mathcal{E}_{1/R}) \leq 4 \sqrt{\frac{(1 - k)(kR^2 + 1)^2}{((1 - k)R^2 + 1)^2}}.$$

Just by inspection, we see that this bound is an increasing function, so we calculate that

$$\begin{aligned} &\leq \lim_{R \rightarrow \infty} 4 \sqrt{\frac{(1 - k)(kR^2 + 1)^2}{((1 - k)R^2 + 1)^2}} \\ &= \frac{4k}{\sqrt{1 - k}} \end{aligned}$$

which shows that for fixed k , $\Lambda \cap \mathcal{E}_{1/R}$ is at most of bounded length, independent of R . Thus, there must exist some constant C_2 such that there are at most C_2 lattice points in $\Lambda \cap \mathcal{E}_{1/R}$. Define $C = \max\{C_1, C_2\}$. We have shown that

$$\max_{a \in \mathbb{R}^2} \#\{B(a, R^{1/3}) \cap \mathcal{E}_{1/R} \cap \mathbb{Z}^2\} < C$$

which proves the first claim.

To prove the second claim of the lemma, we can simply partition $\mathcal{E}_{1/R}$ into sections of length $R^{1/3}$. There will be $R^{2/3}$ of these sections. Applying the previous result,

$$\#\{\mathcal{E}_{1/R} \cap \mathbb{Z}^2\} < R^{2/3} \max_{a \in \mathbb{R}^2} \#\{B(a, R^{1/3}) \cap \mathcal{E}_{1/R} \cap \mathbb{Z}^2\} < CR^{2/3}$$

which finishes the proof of the Lemma.

Remark 5.2. The requirement that the length of the minor and major axes both grow as $R \rightarrow \infty$ is necessary. For example, consider the ellipse centered at the origin with minor axis length 1 and major axis length R . As $R \rightarrow \infty$, the $1/R$ -neighborhood of this ellipse will contain arbitrarily many lattice points on the line $y = 1$.

Remark 5.3. Another fact that's important to remember is that bounded area does not imply bounded number of lattice points. For example, consider the rectangle centered at the origin with side lengths n and $1/n$. The area of the rectangle is uniformly bounded in n ; however, as $n \rightarrow \infty$ it will contain an unbounded number of lattice points.

5.2. A neighborhood of a closed, bounded, convex curve.

Here, we focus on a Proposition of Stein and Wainger [SW99] that was later improved by Mirek in [MST19, MSZK20a, MSZK20b]. There are multiple versions of this proposition presented in Mirek's various papers. In contrast to the previous results we have discussed, this Proposition is concerned with larger neighborhoods.

The original proposition is due to Stein and Wainger.

Proposition 5.4 ([SW99]). *Let $\Omega \subset \mathbb{R}^k$ be a convex set contained in a ball of radius r . Let $N_\Omega = \#\{x \in \Omega \cap \mathbb{Z}^k \mid d(x, \partial\Omega) < s\}$. Then, for $s \geq 1$, $N_\Omega = \mathcal{O}(r^{k-1/2}s^{1/2})$.*

In this proposition, they conjecture that they expect the bound can be improved to $\mathcal{O}(r^{k-1}s)$. In Mirek's sequence of papers, he proves the following versions of the above proposition. Additionally, it is in [MST19] that Mirek suggests there is an error in the proposition of Stein and Wainger. Although this first proposition is weaker than that of [SW99], it is sufficient for their purposes.

We note that in Mirek's publications, he uses \lesssim where we have used \ll .

Proposition 5.5 ([MST19]). *Fix $0 \leq \sigma \leq 1/3$. Let $\Omega \subset \mathbb{R}^k$ be a convex set contained in a ball of radius $r \geq 1$. Let $N_\Omega = \#\{x \in \Omega \cap \mathbb{Z}^k \mid d(x, \partial\Omega) < s\}$. Then,*

- (1) *If $1 \leq s \leq r^{1-3\sigma}$, $N_\Omega = \mathcal{O}(r^{k-\sigma})$.*
- (2) *If $1 \leq s \leq r$ and $\exists x'_0 \in \mathbb{R}^k$ and $c > 0$ such that $\Omega \supseteq B(x'_0, cr)$, $N_\Omega = \mathcal{O}(sr^{k-1})$.*

Proposition 5.6 ([MSZK20a]).

Let $\Omega \subset \mathbb{R}^k$ be a bounded and convex set and let $0 < s \ll \text{diam}(\Omega)$. Then,

$$|\{x \in \mathbb{R}^k \mid d(x, \partial\Omega) < s\}| \ll_k s \text{diam}(\Omega)^{k-1}.$$

Proposition 5.7 ([MSZK20b]).

Let $\Omega \subset \mathbb{R}^k$ be a bounded and convex set and let $1 \leq s \leq \text{diam}(\Omega)$. Then,

$$|\{x \in \mathbb{Z}^k \mid d(x, \partial\Omega) < s\}| \ll_k s \text{diam}(\Omega)^{k-1}.$$

Remark 5.8. In Propositions 5.6, 5.7, the implicit constant depends only on k and not on the convex set Ω .

Since the result of Proposition 5.7 is stronger than that of Proposition 5.5, we will only go through the proofs of Proposition 5.6 and 5.7. If we assume the truth of Proposition 5.6, then Proposition 5.7 is not too hard to prove. Let $\Omega(s) = \{x \in \mathbb{Z}^k \mid d(x, \partial\Omega) < s\}$. Let μ_k denote the Lebesgue measure on \mathbb{R}^k . Then,

$$\#\Omega(s) = \sum_{x \in \Omega(s)} 1 \ll_k \sum_{x \in \Omega(s)} \mu_k(B(x, 1/2))$$

since the $B(x, 1/2)$ are disjoint from the others. For a countable set S we use $\#S$ to denote the number of elements in the set (this may be infinite). The implicit constant in this inequality depends exponentially on k ; for arbitrarily large dimensions, the volume of the sphere becomes arbitrarily small. We can make this constant more explicit by noticing that the formula for the volume of an k -dimensional sphere is

$$\mu_k(B(x, R)) = \frac{\pi^{k/2}}{\Gamma(k/2 + 1)} R^k$$

where Γ denotes the Euler-Gamma function and can be understood as a generalization of the factorial. This shows that the implicit constant C_k is at least

$$C_k \geq \left(\frac{2}{\pi}\right)^{k/2} \Gamma(k/2 + 1).$$

To finish the proof, notice that

$$\bigcup_{x \in \Omega(s)} B(x, 1/2) \subset \{x \in \mathbb{R}^k \mid d(x, \partial\Omega) < s + 1/2\}.$$

Thus, by Proposition 5.6, we know this set has measure $\ll_k (s + 1/2)\text{diam}(\Omega)^{k-1}$. However, we can also write that the measure of the set is $\ll_k s \text{diam}(\Omega)^{k-1}$ for a different

implicit constant. This step is possible because we have assumed that $s \geq 1$; in fact, it suffices to assume $s \gg_k 1/2$ so the Proposition can be strengthened a bit.

To prove Proposition 5.6, we will use a tool from measure theory: the Vitali covering lemma.

Lemma 5.9 (Vitali Covering Lemma). *Let $\{B_k \mid k \in K\}$ be an arbitrary collection of balls in \mathbb{R}^d such that $\sup_{k \in K} \{\text{rad}(B_k)\} < \infty$. Then, there exists a countable subcollection $\{B_k \mid k \in K'\}$ of disjoint balls that satisfies*

$$\bigcup_{k \in K} B_k \subseteq \bigcup_{k \in K'} 5B_k$$

where $5B_k$ denotes the ball with the same center but five times the radius.

Continuing with the proof of Mirek's result, we first let $r = \text{diam}(\Omega)$. Since the size of the neighborhood is not changed by translation, we can assume that $\Omega \subseteq B(0, r)$. Notice that for any $s > 0$,

$$\{x \in \mathbb{R}^k \mid d(x, \partial\Omega) < s\} \subseteq \bigcup_{y \in \partial\Omega} B(y, s).$$

The union on the RHS is uncountable; however, the Vitali covering lemma also applies to uncountable collections of balls. By Vitali there exists some countable collection $Y \subset \partial\Omega$ such that

$$\bigcup_{y \in \partial\Omega} B(y, s) \subseteq \bigcup_{y \in Y} B(y, 5s) \implies \mu_k \left(\bigcup_{y \in \partial\Omega} B(y, s) \right) \leq \mu_k \left(\bigcup_{y \in Y} B(y, 5s) \right).$$

In fact, Y must be *finite* (which is also countable). This is because all of the balls in our collection are of the same radius $s > 0$, so the process of finding a maximal disjoint subcollection will terminate after some finite number of steps. Technically, this proof does not require the Vitali covering lemma in its full strength since we are working with a very "nice" collection of balls. We could directly prove the existence of the finite subcollection Y by using the compactness of $\partial\Omega$ to get a finite subcovering, and then show that once we increase the radii of all these balls by a factor of 5, any point in the s -neighborhood of the boundary will be contained in at least one ball of radius $5s$.

For every $y \in Y$, Mirek then defines $z(y)$ as the nearest point to y on $\partial B(0, r)$. This point is well-defined and unique because we are working with the Euclidean metric. This argument deserves a more detailed explanation; however, due to time constraints, the author will not be able to add these details until a later revision. For an intuitive and geometric point of view, see Figure 3 for an example of how things can go wrong if we use a different metric. Furthermore, since $\Omega \subseteq B(0, r)$, we know that for any $y \neq y' \in Y$, $\|y - y'\| \leq \|z(y) - z(y')\|$ which tells us that the balls $B(z(y), 5s)$ are also disjoint. Importantly, this implies that

$$\mu_k \left(\bigcup_{y \in Y} B(y, 5s) \right) = \mu_k \left(\bigcup_{y \in Y} B(z(y), 5s) \right).$$

Now, we notice that all of these balls centered at the $z(y)$ are contained in the annulus

$$\{x \in \mathbb{R}^k \mid r - 5s < |x| < r + 5s\} =: \text{ann}(0; r - 5s, r + 5s).$$

We can quickly calculate that

$$\mu_k(\text{ann}(0; r - 5s, r + 5s)) = \frac{\pi^{k/2}}{\Gamma(k/2 + 1)} \left((r + 5s)^k - (r - 5s)^k \right)$$

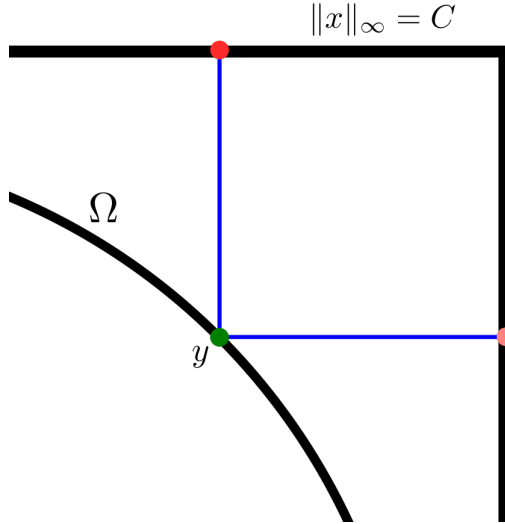


FIGURE 3. The “ball” when using the ∞ -norm is actually a square and this allows for the possibility of a non-unique nearest-point-projection.

$$\begin{aligned}
 &= \frac{\pi^{k/2}}{\Gamma(k/2 + 1)} (10k sr^{k-1} + \mathcal{O}(r^{k-3})) \\
 &\ll_k sr^{k-1}.
 \end{aligned}$$

Since $r = \text{diam}(\Omega)$ this finishes the proof.

6. APPLICATION TO PDES

Definition 6.1. A problem is called “well-posed” if

- (1) A solution exists
- (2) The solution is unique
- (3) The solution’s behavior changes continuously with respect to the initial conditions

If one or more of the criteria is not met, the problem may be called “ill-posed”.

We have glossed over a few details. Firstly, it is also important to distinguish between two types of problems: the Dirichlet problem and the Cauchy problem. The Dirichlet problem is as follows: given a PDE, a region in $S \subset \mathbb{R}^n$, and a function f that is defined on ∂S , is there a unique function u that satisfies the PDE inside s and is equivalent to f on the boundary? We are implicitly assuming some properties of ∂S , essentially that it is sufficiently smooth. On the other hand, the Cauchy problem asks the following: given a PDE and initial condition $u(0, x) = v(x)$, is there a function u that satisfies the PDE subject to the initial condition? Although these two problems seem similar, they lend themselves to different interpretations; there are often settings where one of the problems is ill-posed but the other is well-posed. For a more natural interpretation, one can think of the Cauchy problem as trying to figure out the trajectory of a particle given the initial position and velocity and the Dirichlet problem as trying to figure out the trajectory given the initial and final states. These two types of problems are also commonly referred to as boundary value problems and initial value problems.

Secondly, and more importantly, what actually counts as a solution to a PDE? There are two broad types of solutions: strong and weak. A strong solution to a PDE is a function that satisfies the PDE and is sufficiently differentiable (*i.e.*, all the derivatives in specified exist). A weak solution is one that satisfies the PDE in some precisely defined

equivalent sense, but itself may not be sufficiently differentiable. This may sound absurd, but we can very cleverly pass the burden of differentiability to another function. Loosely speaking, consider the formula for integration by parts

$$\int f'g = fg - \int fg'.$$

We suppose that g is some well-behaved function. If f were differentiable, then the above equation is satisfied. However, if we have two functions f and f' that satisfy the above equation for all g , then we can call f' the “weak” derivative of f . There are many details missing here, but this is the general idea behind weak solutions to a PDE. More generally, when studying PDEs, people are interested in finding weak (or generalized) solutions within certain classes of functions. Although it would be nice to find smooth solutions, it is often more reasonable to look for generalized solutions, and then handle separately the issue of regularity of the weak solutions.

Definition 6.2. The linear (or free) Schrödinger equation is

$$iu_t + \Delta u = 0.$$

The nonlinear (or semilinear) Schrödinger equation (NLSE) is

$$iu_t + \Delta u = |u|^{p-1}u$$

for $p > 1$ and $u(t, x) : \mathbb{R} \times M \rightarrow \mathbb{C}$. M is usually Euclidean space or a Torus. If M is a Torus, this problem is called the *periodic* NLSE (since u will be periodic in each component of x).

Here, $\Delta u = \sum_{j=1}^n u_{x_j x_j}$ denotes the Laplacian operator. In other words, Δu is the sum of all the non-mixed, second order partial derivatives of u .

The problem Bourgain and many others were interested in is the following initial value problem regarding the well-posedness of different forms of the NLSE.

- (Existence) Given suitable initial data, does there exist a solution to the initial value problem that exists for some time period $[0, T]$, $T > 0$? What can we say about T ? How well-behaved (regular) is the solution?
- (Uniqueness) Is it possible to have two solutions to the same initial value problem on the same time interval in the same class of functions? To what extent does the assumed regularity of the solutions affect the answer to the previous question?
- (Continuous dependence on data) If we slightly perturb the initial data, what happens to the solution on the time interval of existence?

In his work involving the periodic nonlinear Schrödinger equation (NLSE) on the Torus, Jean Bourgain was the first to use tools from analytic number theory to prove local well-posedness of some specific forms of the NLS. Bourgain and others were able to reduce the problem of proving well-posedness to the problem of bounding the number of lattice points on convex curves. The connection between these two topics is not intuitive at all! Furthermore, there are multiple different approaches to proving local well-posedness. In the specific case that’s relevant to us, bounds on lattice points are needed for a type of inequality called a Strichartz estimate. These types of estimates are bounds on the norm of solutions to the NLSE, and they are useful for analyzing the fixed-point problem that arises from transforming the NLSE into an integral equation. Essentially, the integral version of the NLSE defines an operator whose fixed point is a solution to the NLSE; to apply the contraction mapping principle, we first need to find a space of functions that is mapped into itself by the operator and this is where Strichartz estimates are useful. There are a lot of details that have been omitted in this section, in fact one could write

an entire report on this topic, but hopefully the information provided is useful for those who are interested in applications of the main topic of this essay.

REFERENCES

- [Bou09] J. Bourgain. *On Strichartz's Inequalities and the Nonlinear Schrödinger Equation on Irrational Tori*, pages 1–20. Princeton University Press, 2009.
- [CC92] J. Cilleruelo and A. Córdoba. Trigonometric polynomials and lattice points. *Proc. Amer. Math. Soc.*, 115(4):899–905, 1992.
- [CG09] Javier Cilleruelo and Andrew Granville. Close lattice points on circles. *Canadian Journal of Mathematics*, 61(6):1214–1238, 2009.
- [Jar26] V. V. Jarnik. Über die gitterpunkte auf konvexen kurven. *Mathematische Zeitschrift*, 24:500–518, 1926.
- [MST19] Mariusz Mirek, Elias M. Stein, and Bartosz Trojan. $\ell^p(F^d)$ -estimates for discrete operators of radon type: Maximal functions and vector-valued estimates. *Journal of Functional Analysis*, 277(8):2471–2521, 2019.
- [MSZK20a] Mariusz Mirek, Elias M. Stein, and Pavel Zorin-Kranich. A bootstrapping approach to jump inequalities and their applications. *Analysis & PDE*, 13(2):527–558, Mar 2020.
- [MSZK20b] Mariusz Mirek, Elias M. Stein, and Pavel Zorin-Kranich. Jump inequalities for translation-invariant operators of radon type on \mathbb{Z}^d . *Advances in Mathematics*, 365, 2020.
- [SW99] Elias M. Stein and Stephen Wainger. Discrete analogues in harmonic analysis, i: ℓ^2 estimates for singular radon transforms. *American Journal of Mathematics*, 121(6):1291–1336, 1999.