

# Chapter 4

## Superirregularity

József Beck

**Abstract** Finding the integer solutions of a Pell equation is equivalent to finding the integer lattice points in a long and narrow tilted hyperbolic region, where the slope is a quadratic irrational. Motivated by this relationship, we carry out here a systematic study of point counting with respect to translated or congruent families of any given long and narrow hyperbolic region. First we discuss the important special case when the underlying point set is the set of integer lattice points in the plane and the slope of the given hyperbolic region is arbitrary but fixed; see Theorems 3–21. Then we switch to the general case of an arbitrary point set of density one in the plane, and study point counting with respect to congruent copies of a given hyperbolic region; see Theorem 30. The main results are about the extra large discrepancy that we call *superirregularity*. This means that there is always a translated/congruent copy of any given long and narrow hyperbolic region of large area, for which the actual number of points in the copy differs from the area *as much as possible*, i.e. the discrepancy is at least a constant multiple of the area. Our theorems demonstrate, in a quantitative sense, that in point counting with respect to translated/congruent copies of any long and narrow hyperbolic region, superirregularity is inevitable.

### 4.1 Introduction

*Notation.* For any real valued function  $f$  and positive function  $g$ , we write  $f = O(g)$  to indicate that there exists a positive constant  $c$  such that  $|f| < cg$ , and also write  $f = o(g)$  to indicate that  $f/g \rightarrow 0$ . We write  $\|z\|$  to denote the distance of a real number  $z$  to the nearest integer. Furthermore,  $c_0, c_1, c_2, \dots$  denote positive constants which may depend on some of the parameters that arise from our discussion.

---

J. Beck (✉)

Department of Mathematics, Rutgers University, New Brunswick, NJ 08903, USA  
e-mail: [jbeck@math.rutgers.edu](mailto:jbeck@math.rutgers.edu)

### 4.1.1 Pell's Equation: Bounded Fluctuations

Our starting point is the well-known Pell's equation, a standard part of any introductory course on number theory. The theory of Pell's equation, while mostly elementary, is nevertheless one of the most beautiful chapters in the whole of mathematics. Also, it is very important, since the concept of *units* plays a key role in algebraic number theory.

We briefly recall the main results. Consider, for simplicity, the concrete equation  $x^2 - 2y^2 = \pm 1$ . This equation has infinitely many integral solutions; in fact, the set of all integral solutions  $(x_k, y_k) \in \mathbf{Z}^2$  forms a cyclic group generated by the least positive solution. More precisely, we have

$$x_k + y_k\sqrt{2} = \pm(1 + \sqrt{2})^k, \quad k \in \mathbf{Z}.$$

All integral solutions of  $x^2 - 2y^2 = 1$  are given by  $x_k + y_k\sqrt{2} = \pm(1 + \sqrt{2})^{2k}$ , while all integral solutions of  $x^2 - 2y^2 = -1$  are given by  $x_k + y_k\sqrt{2} = \pm(1 + \sqrt{2})^{2k+1}$ . In particular, all positive integer solutions of  $x^2 - 2y^2 = 1$  are given by

$$x_k + y_k\sqrt{2} = (1 + \sqrt{2})^{2k} = (3 + 2\sqrt{2})^k, \quad k = 1, 2, 3, \dots$$

Taking the algebraic conjugate  $x_k - y_k\sqrt{2} = (3 - 2\sqrt{2})^k$ , and combining these two equations, we obtain the explicit formulas

$$x_k = \frac{(3 + 2\sqrt{2})^k + (3 - 2\sqrt{2})^k}{2} \quad \text{and} \quad y_k = \frac{(3 + 2\sqrt{2})^k - (3 - 2\sqrt{2})^k}{2\sqrt{2}}.$$

Since  $0 < 3 - 2\sqrt{2} < \frac{1}{5}$ , we have

$$x_k = \text{the nearest integer to } \frac{1}{2}(3 + 2\sqrt{2})^k$$

and

$$y_k = \text{the nearest integer to } \frac{1}{2\sqrt{2}}(3 + 2\sqrt{2})^k.$$

If  $k$  is large, the error is very small. For example, the 10-th solution of  $x^2 - 2y^2 = 1$  in positive integers is the pair  $x_{10} = 22,619,537$  and  $y_{10} = 15,994,428$ . Here we find

$$\frac{1}{2}(3 + 2\sqrt{2})^{10} = 22619536.99999998895 \dots$$

and

$$\frac{1}{2\sqrt{2}}(3 + 2\sqrt{2})^{10} = 15994428.000000007815 \dots$$

Let  $F(N) = F(\sqrt{2}; 1; N)$  denote the number of positive integer solutions of the Pell equation  $x^2 - 2y^2 = 1$  up to  $N$ , in the sense<sup>1</sup> that  $x \geq 1$  and  $1 \leq y \leq N$ . We have

$$k \leq F(N) \quad \text{if and only if} \quad \frac{(3 + 2\sqrt{2})^k - (3 - 2\sqrt{2})^k}{2\sqrt{2}} \leq N,$$

which implies the asymptotic formula

$$F(N) = F(\sqrt{2}; 1; N) = \frac{\log N}{\log(3 + 2\sqrt{2})} + O(1). \tag{4.1}$$

The formula (4.1) says that the counting function  $F(N) = F(\sqrt{2}; 1; N)$  has an extremely predictable, almost deterministic behavior: it is  $c_2 \log N$  plus some bounded error term.

Note that (4.1) has some far-reaching generalizations. Let  $[\gamma_1, \gamma_2]$  be an arbitrary interval, and let  $F(\sqrt{2}; [\gamma_1, \gamma_2]; N)$  denote the number of positive integer solutions of the Pell inequality  $\gamma_1 \leq x^2 - 2y^2 \leq \gamma_2$ , with  $x \geq 1$  and  $1 \leq y \leq N$ . By using the theory of indefinite binary quadratic forms, it is easy to prove the following analog of (4.1). We have

$$F(\sqrt{2}; [\gamma_1, \gamma_2]; N) = c_0 \log N + O(1), \tag{4.2}$$

where the constant factor  $c_0 = c_0(\sqrt{2}; \gamma_1, \gamma_2)$  is independent of  $N$ .

Furthermore, we can switch from  $\sqrt{2}$  to any other *quadratic irrational*  $\alpha$ . This means that  $\alpha$  is a root of a quadratic equation  $Ax^2 + Bx + C = 0$  with integral coefficients such that the discriminant  $B^2 - 4AC \geq 2$  is not a complete square. An equivalent definition is that  $\alpha = (a + \sqrt{d})/b$  for some integers  $a, b, d$  such that  $b \neq 0$  and  $d \geq 2$  is not a complete square. Note that the quadratic irrationals are characterized by their continued fractions. The continued fractions of  $\alpha$  is finally periodic if and only if  $\alpha$  is a quadratic irrational. For example,

$$\frac{24 - \sqrt{15}}{17} = 1 + \frac{1}{5 + \frac{1}{2 + \frac{1}{3 + \frac{1}{2 + \frac{1}{3 + \dots}}}}} = [1; 5, 2, 3, 2, 3, 2, 3, \dots] = [1; 5, \overline{2, 3}].$$

Let us go back to (4.2) and to the special case  $\alpha = \sqrt{2}$ . If  $-2 < \gamma_1 \leq -1$  and  $1 \leq \gamma_2 < 2$ , then

$$c_0(\sqrt{2}; \gamma_1, \gamma_2) = \frac{1}{\log(1 + \sqrt{2})} = \frac{2}{\log(3 + 2\sqrt{2})}. \tag{4.3}$$

---

<sup>1</sup>For simplicity of notation, it is more convenient to restrict the second variable  $y$ .

If  $-1 < \gamma_1 \leq 1 \leq \gamma_2 < 2$ , then

$$c_0(\sqrt{2}; \gamma_1, \gamma_2) = \frac{1}{\log(3 + 2\sqrt{2})}. \quad (4.4)$$

Finally, if  $-1 < \gamma_1 \leq \gamma_2 < 1$ , then of course

$$c_0(\sqrt{2}; \gamma_1, \gamma_2) = 0. \quad (4.5)$$

### 4.1.2 The Naive Area Principle

It is very interesting to compare these well-known asymptotic results about the number of solutions of the Pell equation/inequality to what we like to call the *Naive Area Principle*, a natural guiding intuition in lattice point theory. It goes roughly as follows. If a nice region has a large area, then it should contain a large number of lattice points, and the number of lattice points is close to the area.

Of course, the heart of the matter is how we define a nice region precisely. Consider, for example, the infinite open horizontal strip of height one, given by  $0 < y < 1$ ,  $-\infty < x < \infty$ . It has infinite area, but it does not contain any lattice point. The reader is likely to agree that the infinite strip is a nice region, so the Naive Area Principle is clearly violated here.

A less trivial example comes from the Pell inequality

$$-\frac{1}{2} \leq x^2 - 2y^2 \leq \frac{1}{2}. \quad (4.6)$$

This is a hyperbolic region of infinite area, and contains no lattice point except the origin. The reader is again likely to agree that the hyperbolic region (4.6) is also nice, so this is again a violation of the Naive Area Principle.

Next we switch from (4.6) to the general Pell inequality

$$\gamma_1 \leq x^2 - 2y^2 \leq \gamma_2, \quad (4.7)$$

where  $-\infty < \gamma_1 < \gamma_2 < \infty$  are arbitrary real numbers. Of course, the hyperbolic region (4.7) has infinite area. What we want to compute is the area of a finite segment. Consider the finite region

$$H(\sqrt{2}; [\gamma_1, \gamma_2]; N) = \{(x, y) \in \mathbf{R}^2 : \gamma_1 \leq x^2 - 2y^2 \leq \gamma_2, x \geq 1, 1 \leq y \leq N\}. \quad (4.8)$$

If  $N$  is very large compared to the pair of constants  $\gamma_1$  and  $\gamma_2$ , then the finite region  $H(\sqrt{2}; [\gamma_1, \gamma_2]; N)$  looks like a *hyperbolic needle*. It is easy to give a good estimate for the area of this hyperbolic needle. We have

$$\text{area}(H(\sqrt{2}; [\gamma_1, \gamma_2]; N)) = \frac{\gamma_2 - \gamma_1}{2\sqrt{2}} \log N + O(1), \quad (4.9)$$

where the implicit constant in the term  $O(1)$  is independent of  $N$ , but may depend on  $\gamma_1$  and  $\gamma_2$ .

The proof of (4.9) is based on the familiar factorization

$$x^2 - 2y^2 = (x + y\sqrt{2})(x - y\sqrt{2}), \quad (4.10)$$

and on the computation of the Jacobian of the corresponding substitution; this explains the factor  $2\sqrt{2}$  in the denominator in (4.9). The details are easy, and go as follows. In view of the factorization (4.10), it is more convenient to compute the area of the following slight variant of the region (4.9). Let

$$\begin{aligned} H^*(\sqrt{2}; [\gamma_1, \gamma_2]; N) \\ = \{(x, y) \in \mathbf{R}^2 : \gamma_1 \leq x^2 - 2y^2 \leq \gamma_2, 1 \leq x + y\sqrt{2} \leq 2\sqrt{2}N\}. \end{aligned} \quad (4.11)$$

Consider the substitution

$$u_1 = x + y\sqrt{2}, \quad u_2 = x - y\sqrt{2}, \quad (4.12)$$

which is equivalent to

$$x = \frac{u_1 + u_2}{2}, \quad y = \frac{u_1 - u_2}{2\sqrt{2}}.$$

The corresponding determinant is

$$\frac{\partial(u, v)}{\partial(x, y)} = \begin{vmatrix} 1 - \sqrt{2} \\ 1 \quad \sqrt{2} \end{vmatrix} = 2\sqrt{2}.$$

Applying the substitution (4.12), we have

$$\begin{aligned} \text{area}(H^*(\sqrt{2}; [\gamma_1, \gamma_2]; N)) &= \frac{1}{2\sqrt{2}} \int_1^{2\sqrt{2}N} \left( \int_{\gamma_1/u_1}^{\gamma_2/u_1} du_2 \right) du_1 \\ &= \frac{1}{2\sqrt{2}} \int_1^{2\sqrt{2}N} \frac{\gamma_2 - \gamma_1}{u_1} du_1 = \frac{\gamma_2 - \gamma_1}{2\sqrt{2}} \log N + O(1). \end{aligned} \quad (4.13)$$

Simple geometric consideration shows that

$$\text{area}(H(\sqrt{2}; [\gamma_1, \gamma_2]; N)) = \text{area}(H^*(\sqrt{2}; [\gamma_1, \gamma_2]; N)) + O(1),$$

and so (4.13) implies (4.9).

Now let us return to the Naive Area Principle. Comparing (4.2), (4.8) and (4.9), it is reasonable to expect, in view of the Naive Area Principle, that the counting function  $F(\sqrt{2}; [\gamma_1, \gamma_2]; N)$  is close to the area of the hyperbolic needle  $H(\sqrt{2}; [\gamma_1, \gamma_2]; N)$ . In other words, it is reasonable to expect that

$$c_0(\sqrt{2}; \gamma_1, \gamma_2) = \frac{\gamma_2 - \gamma_1}{2\sqrt{2}}. \tag{4.14}$$

Unfortunately, the Naive Area Principle is almost always violated in the quantitative sense that (4.14) fails for the overwhelming majority of the choices  $-\infty < \gamma_1 < \gamma_2 < \infty$ . In fact, the two sides of (4.14) have completely different behavior. The left-hand side has discrete jumps and the right-hand side is a continuous function of  $\gamma_1$  and  $\gamma_2$ . For example, as  $\gamma_1$  and  $\gamma_2$  run in the interval  $-2 < \gamma_1 < \gamma_2 < 2$ , the constant factor  $c_0(\sqrt{2}; \gamma_1, \gamma_2)$  has only 3 possible values, namely

$$0, \quad \frac{1}{\log(3 + 2\sqrt{2})}, \quad \frac{2}{\log(3 + 2\sqrt{2})};$$

see (4.3)–(4.5). This shows, in a quantitative way, how the general Pell inequality (4.7) violates the Naive Area Principle.

### 4.1.3 The Giant Leap in the Inhomogeneous Case: Extra Large Fluctuations

Using the familiar factorization (4.10), we can rewrite the Pell equation  $x^2 - 2y^2 = \pm 1$ , restricted to positive integers, as

$$|x^2 - 2y^2| \leq 1 \quad \text{or} \quad |y\sqrt{2} - x|(y\sqrt{2} + x) \leq 1 \quad \text{or} \quad \|y\sqrt{2}\|(y\sqrt{2} + x) \leq 1, \tag{4.15}$$

where  $\|z\|$  denotes, as usual, the distance of a real number  $z$  from the nearest integer. Notice that in (4.15),  $x$  is the nearest integer to  $y\sqrt{2}$ , which is an irrational number. Since  $y\sqrt{2} \approx x$ , the inequality (4.15) is basically equivalent to the vague inequality

$$\|y\sqrt{2}\| \leq \frac{1 + o(1)}{2\sqrt{2}y}. \tag{4.16}$$

The vagueness of (4.16) comes from the additional term  $o(1)$ , which tends to 0 as  $y \rightarrow \infty$ . The formula (4.16) is ambiguous, but surely every mathematician understands what we are talking about here.

An expert in number theory would classify (4.16) as a typical problem in diophantine approximation. Next we give a nutshell summary of diophantine approximation.

The classical problem in the theory of diophantine approximation is to find good rational approximations of irrational numbers. More precisely, we want to decide whether an inequality

$$\|n\alpha\| < \frac{1}{n\varphi(n)} \quad \text{or} \quad \left| \alpha - \frac{m}{n} \right| < \frac{1}{n^2\varphi(n)}, \quad (4.17)$$

or in general,

$$\|n\alpha - \beta\| < \frac{1}{n\varphi(n)}, \quad (4.18)$$

where  $\alpha$  is a given irrational number and  $\beta$  is a given real number, has infinitely many integral solutions in  $n$ , and if this is the case, to determine the solutions, or at least the asymptotic number of integral solutions. Here  $\varphi(n)$  is a positive increasing function of  $n$ .

The diophantine inequality (4.17) is said to be homogeneous, whereas the diophantine inequality (4.18) is said to be inhomogeneous. For example, in the homogeneous case, the best possible result is Hurwitz's well-known theorem, that for any irrational number  $\alpha$ , the inequality

$$\|n\alpha\| < \frac{1}{\sqrt{5}n}$$

has infinitely many positive integer solutions.

In the inhomogeneous case, we can mention an old result of Kronecker, that for any irrational number  $\alpha$  and any real number  $\beta$ , the inequality

$$\|n\alpha - \beta\| < \frac{3}{n}$$

has infinitely many positive integer solutions. Perhaps the strongest inhomogeneous result is Minkowski's theorem, that for any irrational number  $\alpha$ , the inequality

$$\|n\alpha - \beta\| < \frac{1}{4n}$$

has infinitely many integer but not necessarily positive solutions, unless  $0 < \beta < 1$  is an integer multiple of  $\alpha$  modulo one.

The homogeneous case (4.17) has a complete theory based on the effectiveness of the tool of continued fractions. These are classical results due mostly to Euler and Lagrange. Unfortunately, we know much less about the inhomogeneous case. Very recently, the author proved some new results in this direction, and basically covered the case when  $\alpha$  is an arbitrary quadratic irrational and  $\beta$  is a typical real number. These results form a large part of the forthcoming book [2]; see also the recent papers [8, 9].

Before formulating our main results, we want to first elaborate on the connection between homogeneous/inhomogeneous diophantine inequalities, such as (4.17) and (4.18), and homogeneous/inhomogeneous Pell inequalities.

#### 4.1.3.1 Homogeneous and Inhomogeneous Pell Inequalities

The general form of a quadratic curve on the plane is

$$a_{11}x^2 + a_{12}xy + a_{22}y^2 + a_{13}x + a_{23}y + a_{33} = 0. \quad (4.19)$$

We are interested in the integral solutions  $(x, y) \in \mathbf{Z}^2$  of an arbitrary inequality

$$\gamma_1 \leq a_{11}x^2 + a_{12}xy + a_{22}y^2 + a_{13}x + a_{23}y \leq \gamma_2, \quad (4.20)$$

where  $\gamma_1 < \gamma_2$  are given real numbers. Note that the inequality (4.20) defines a plane region, and the boundary consists of two curves of the type (4.19). In the case of negative discriminant  $D = a_{12}^2 - 4a_{11}a_{22} < 0$ , the inequality (4.20) defines a bounded region where the boundary curves are two ellipses. The case of positive discriminant  $D = a_{12}^2 - 4a_{11}a_{22} > 0$  is much more interesting, because then the inequality (4.20) defines an unbounded region, where the boundary curves are two hyperbolas, and thus we have a chance for infinitely many integral solutions of (4.20).

For simplicity, assume that the coefficients  $a_{11}, a_{12}, a_{22}$  in (4.20) are integers and  $D = a_{12}^2 - 4a_{11}a_{22} > 0$ . We can factorize the quadratic part in the form

$$a_{11}x^2 + a_{12}xy + a_{22}y^2 = a_{11}(x - \alpha y)(x - \alpha' y), \quad (4.21)$$

where

$$\alpha = \frac{-a_{12} + \sqrt{D}}{2a_{11}} \quad \text{and} \quad \alpha' = \frac{-a_{12} - \sqrt{D}}{2a_{11}}. \quad (4.22)$$

Using (4.21), we can rewrite (4.20) in the form

$$\gamma_1 \leq (x - \alpha y + \rho_1)(x - \alpha' y + \rho_2) \leq \gamma_2, \quad (4.23)$$

where

$$\rho_1 + \rho_2 = \frac{a_{13}}{a_{11}} \quad \text{and} \quad \alpha' \rho_1 + \alpha \rho_2 = -\frac{a_{23}}{a_{11}}.$$

Note that  $\gamma_1, \gamma_2$  are generic numbers; the pair  $\gamma_1, \gamma_2$  in (4.20) is not necessarily the same as the pair  $\gamma_1, \gamma_2$  in (4.23).



Without loss of generality we can assume<sup>2</sup> that  $|a_{12}| \leq a_{11} \leq \sqrt{D/3}$ , and then we have  $\alpha > 0 > \alpha'$ .

For simplicity, assume that the interval  $[\gamma_1, \gamma_2]$  is symmetric with respect to 0, so that it is of the form  $[\gamma_1, \gamma_2] = [-\gamma, \gamma]$ . Assume also that we are interested in the positive integral solutions of (4.23). Since  $\alpha > 0 > \alpha'$ , for large positive  $x$  and  $y$ , the second factor  $(x - \alpha'y + \rho_2)$  in (4.23) is also large and positive, implying that the first factor  $(x - \alpha y + \rho_1)$  in (4.23) has to be very small. In other words,  $x$  has to be the nearest integer to  $(\alpha y - \rho_1)$ . It follows that the symmetric version of (4.20), namely

$$-\gamma \leq a_{11}x^2 + a_{12}xy + a_{22}y^2 + a_{13}x + a_{23}y \leq \gamma,$$

where  $\gamma > 0$  is a given real number, is equivalent to the diophantine inequality

$$\|y\alpha - \rho_1\| < \frac{c}{y + O(1)}, \quad \text{where } c = \frac{\gamma}{\alpha - \alpha'} = \frac{\gamma a_{11}}{\sqrt{D}}. \quad (4.24)$$

Let us return to the inequality (4.20). If the linear part  $a_{13}x + a_{23}y$  in the middle is missing, i.e.  $a_{13} = a_{23} = 0$ , then we have a complete theory based on Pell's equation. More precisely, write  $Q(x, y) = a_{11}x^2 + a_{12}xy + a_{22}y^2$ . Then  $\gamma_1 \leq Q(x, y) \leq \gamma_2$  if and only if

$$Q(x, y) = m \quad \text{for some } m \in \mathbf{Z} \text{ satisfying } \gamma_1 \leq m \leq \gamma_2.$$

We have a complete characterization of the integral solutions of  $Q(x, y) = m$  for any integer  $m$  as follows. For any integer  $m$ , there is a finite list of primary solutions, say,  $(x_j, y_j)$ ,  $j \in J$ , where  $|J| < \infty$ , such that every solution  $x = u$ ,  $y = v$  of  $Q(x, y) = m$  can be written in the form

$$u - \alpha v = \pm \left( \frac{u_0 + v_0 \sqrt{D}}{2} \right)^n (x_j - \alpha y_j)$$

for some  $j \in J$  and  $n \in \mathbf{Z}$ , where  $x = u_0 > 0$ ,  $y = v_0 > 0$  is the least positive solution of Pell's equation  $x^2 - Dy^2 = 4$ . As a byproduct, we deduce<sup>3</sup> that the number of positive integral solutions of the inequality

$$\gamma_1 \leq Q(x, y) \leq \gamma_2, \quad 1 \leq x \leq N, \quad 1 \leq y \leq N$$

has the simple asymptotic form  $c \log N + O(1)$ , where  $c = c(a_{11}, a_{12}, a_{22}, \gamma_1, \gamma_2)$  is a constant and the error term  $O(1)$  is uniformly bounded as  $N \rightarrow \infty$ .

<sup>2</sup>This is a well-known fact from the reduction theory of binary quadratic forms. We omit the proof; see, for example, [31].

<sup>3</sup>For a more detailed proof; see [23].

Exactly the same holds if there is a non-zero linear part  $a_{13}x + a_{23}y$  in (4.20), but its effect cancels out. Note that  $\rho_1$  in (4.23) is an integer.

Finally, if  $\rho_1$  is not an integer, then we say that (4.23) is an inhomogeneous Pell inequality. In view of (4.24), an inhomogeneous Pell inequality (4.23) is basically equivalent to an inhomogeneous diophantine inequality

$$\|n\alpha - \beta\| < \frac{c}{n} \tag{4.25}$$

with  $c = \gamma a_{11}/\sqrt{D}$ , where  $\alpha$  is a quadratic irrational defined in (4.22). The inequality (4.25) is a special case of (4.18) where  $\varphi(n)$  is a constant.

### 4.1.3.2 Some Results

One of the main results in the forthcoming book [2] describes the asymptotic behavior of the number of positive integral solutions of (4.20) for every non-square integer discriminant  $D > 0$  and almost all  $a_{13}, a_{23}$ . The number of solutions

- exhibits extra large fluctuations, proportional to the area,
- satisfies an elegant Central Limit Theorem, and
- satisfies a shockingly precise Law of the Iterated Logarithm; see Theorems 3, A and B below.

For notational simplicity, we formulate the results in the special case of discriminant  $D = 8$ , which corresponds to the most famous quadratic irrational  $\alpha = \sqrt{2}$ .

Since the class number of the discriminant  $D = 8$  is one, the general form of an inhomogeneous Pell inequality of discriminant  $D = 8$  is

$$\gamma_1 \leq (x + \beta_1)^2 - 2(y + \beta_2)^2 \leq \gamma_2, \tag{4.26}$$

where  $\gamma_1 < \gamma_2$  and  $\beta_1, \beta_2 \in [0, 1)$  are fixed constants. For notational simplicity, we restrict ourselves to symmetric intervals  $[-\gamma, \gamma]$  in (4.26); note that everything works similarly for general intervals  $[\gamma_1, \gamma_2]$ .

The factorization

$$(x + \beta_1)^2 - 2(y + \beta_2)^2 = (x + \beta - y\sqrt{2})(x + \beta' + y\sqrt{2}), \tag{4.27}$$

where  $\beta = \beta_1 - \beta_2\sqrt{2}$  and  $\beta' = \beta_1 + \beta_2\sqrt{2}$ , clearly indicates that the asymptotic number of integral solutions of (4.26) depends heavily on the local behavior of  $n\sqrt{2} \pmod 1$ . In fact, (4.26) is essentially equivalent to the inhomogeneous diophantine inequality

$$\|n\sqrt{2} - \beta\| < \frac{c}{n}, \tag{4.28}$$

with  $c = \gamma/2\sqrt{2}$ .

To turn the vague term *essentially equivalent* into a precise statement, we proceed as follows. Let  $F(\sqrt{2}; \beta_1, \beta_2; \gamma; N)$  be the number of integral solutions  $(x, y) \in \mathbf{Z}^2$  of the inequality (4.26) with  $\gamma_2 = \gamma$  and  $\gamma_1 = -\gamma$  satisfying  $1 \leq y \leq N$  and  $x \geq 1$ . It means counting lattice points in a long and narrow hyperbola segment. Next let  $f(\sqrt{2}; \beta; c; N)$  denote the number of integral solutions  $n$  of the inequality (4.28) satisfying  $1 \leq n \leq N$ , where  $\beta = \beta_1 - \beta_2\sqrt{2}$ . Now *essentially equivalent* means that for almost all pairs  $\beta_1, \beta_2$ , we have  $F(\sqrt{2}; \beta_1, \beta_2; \gamma; N) - f(\sqrt{2}; \beta; c; N) = O(1)$  as  $N \rightarrow \infty$ , where  $c = \gamma/2\sqrt{2}$ . More precisely, we have

**Lemma 1.** *Let  $\gamma > 0$  and  $\beta_2$  be arbitrary real numbers. Then for almost all  $\beta_1$ , there exists a finite  $0 < C(\beta_1, \beta_2, \gamma) < \infty$  such that*

$$\int_0^1 C(\beta_1, \beta_2, \gamma) d\beta < \infty$$

and

$$|F(\sqrt{2}; \beta_1, \beta_2; \gamma; N) - f(\sqrt{2}; \beta; c; N)| < C(\beta_1, \beta_2, \gamma)$$

for all  $N \geq 1$ , where  $c = \gamma/2\sqrt{2}$  and  $\beta = \beta_1 - \beta_2\sqrt{2}$ .

We postpone the simple proof to Sect. 4.3.

In view of Lemma 1, it suffices to study the special case  $\beta_2 = 0$  and  $\beta_1 = \beta$ . We have

$$-\gamma \leq (x + \beta)^2 - 2y^2 \leq \gamma, \tag{4.29}$$

where  $\gamma > 0$  and  $\beta \in [0, 1)$  are fixed constants. For simplicity, let  $F(\sqrt{2}; \beta; \gamma; N)$  denote the number of integral solutions  $(x, y) \in \mathbf{Z}^2$  of (4.29) satisfying  $1 \leq y \leq N$  and  $x \geq 1$ . Note that  $F(\sqrt{2}; \beta; \gamma; N)$  counts the number of lattice points in a long and narrow hyperbola segment, or hyperbolic needle, located along a line<sup>4</sup> of slope  $1/\sqrt{2}$ ; see Fig. 4.1.

In the special case  $\gamma = 1$  and  $\beta = 0$ , the inequality (4.29) becomes the simplest Pell equation  $x^2 - 2y^2 = \pm 1$ . The integral solutions  $(x_k, y_k)$  form a cyclic group generated by the smallest positive solution  $x = y = 1$  in the well-known way. We have  $x_k + y_k\sqrt{2} = (1 + \sqrt{2})^k$ , implying the familiar asymptotic formula

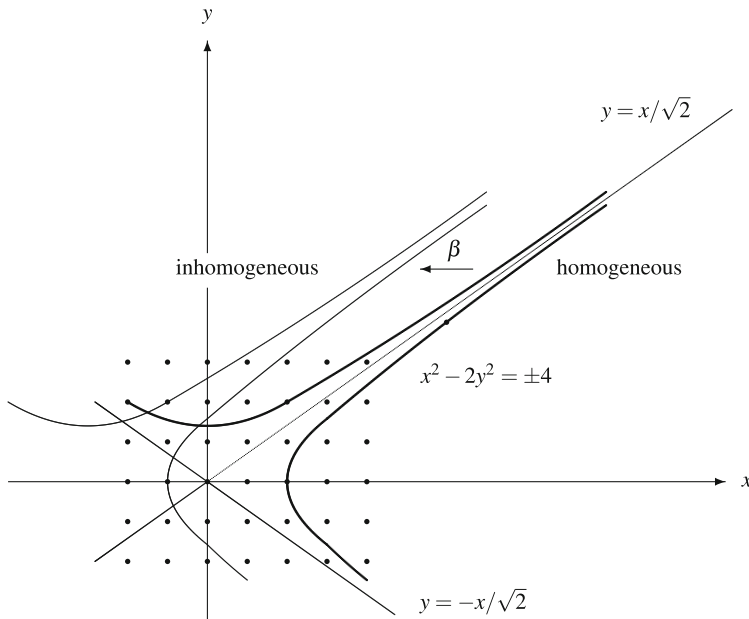
$$F(\sqrt{2}; \beta = 0; \gamma = 1; N) = \frac{\log N}{\log(1 + \sqrt{2})} + O(1), \tag{4.30}$$

where  $1 + \sqrt{2}$  is the fundamental unit of the real quadratic field  $\mathbf{Q}(\sqrt{2})$ .

In sharp contrast to the bounded fluctuation in the homogeneous case  $\beta = 0$ , the inhomogeneous case can exhibit extra large fluctuations proportional to the area;

---

<sup>4</sup>If  $\beta = 0$ , then the line is  $y = x/\sqrt{2}$ .



**Fig. 4.1** A hyperbolic needle

see Theorem 3 below. To explain this, first we have to compute the mean value of  $F(\sqrt{2}; \beta; \gamma; N)$  as  $\beta$  runs through the unit interval  $0 \leq \beta < 1$ .

**Lemma 2.** *We have*

$$\int_0^1 F(\sqrt{2}; \beta; \gamma; N) \, d\beta = \frac{\gamma}{\sqrt{2}} \log N + O(1), \tag{4.31}$$

where the implicit constant in the term  $O(1)$  is independent of  $N$ , but may depend on  $\gamma$ . Moreover, for an arbitrary subinterval  $0 \leq a < b \leq 1$ , we have

$$\lim_{N \rightarrow \infty} \frac{\frac{1}{b-a} \int_a^b F(\sqrt{2}; \beta; \gamma; N) \, d\beta}{\log N} = \frac{\gamma}{\sqrt{2}}. \tag{4.32}$$

The estimates (4.31) and (4.32) express the almost trivial geometric fact that the average number of lattice points contained in all the translated copies of a given region, a hyperbola segment in our special case, is precisely the area of the region; see Lemma 5. We shall give a detailed proof of Lemma 2 in Sect. 4.3.

Now we are ready to formulate our first, and weakest, extra large fluctuation result, demonstrating that the fluctuations can be proportional to the area. This result is hardly more than a warmup for, or simplest illustration of, the main results that will come later.

**Theorem 3.** For  $\gamma = \frac{1}{2}$ , there are continuum many divergence points  $\beta^* \in [0, 1)$  in the sense that

$$\limsup_{n \rightarrow \infty} \frac{F(\sqrt{2}; \beta^*; \gamma = 1/2; n)}{\log n} > \liminf_{n \rightarrow \infty} \frac{F(\sqrt{2}; \beta^*; \gamma = 1/2; n)}{\log n}. \tag{4.33}$$

Note that the fluctuation  $c_3 \log n$  in  $F(\sqrt{2}; \beta^*; \gamma = 1/2; n)$  is as large as possible, apart from a constant factor. This follows from Lemma 4 in the next section. It is fair to say that Theorem 3 represents a sophisticated violation of the Naive Area Principle.

We postpone the proof of Theorem 3 to Sect. 4.3.

Note that Theorem 3 has a far-reaching generalization. It holds for every  $\gamma > 0$ , and we actually have the stronger inequality

$$\limsup_{n \rightarrow \infty} \frac{F(\sqrt{2}; \beta^*; \gamma; n)}{\log n} > \frac{\gamma}{\sqrt{2}} > \liminf_{n \rightarrow \infty} \frac{F(\sqrt{2}; \beta^*; \gamma; n)}{\log n}. \tag{4.34}$$

We shall return to this in Sect. 4.4; see Theorem 12.

Another far-reaching generalization of Theorem 3 will be discussed in Sect. 4.9; see Theorem 21.

Finally, an extra large fluctuation type result for arbitrary point sets, instead of the set  $\mathbf{Z}^2$  of lattice points, will be discussed in Sect. 4.10; see Theorem 30.

We refer to these extra large fluctuation type results as *superirregularity*.

## 4.2 Defending the Naive Area Principle

The estimate (4.30) and inequality (4.33) display the two extreme cases: (1) the negligible bounded fluctuations around the main value which is a constant multiple of  $\log N$ ; and (2) the extra large fluctuations proportional to the area. But what kind of fluctuations do we have for a typical  $\beta$  satisfying  $0 < \beta < 1$ ? We show that for a typical  $\beta$ , the asymptotic number of solutions  $F(\sqrt{2}; \beta; \gamma; N)$ , as  $N \rightarrow \infty$ , justifies the Naive Area Principle. And beyond that, a more thorough look reveals randomness.

Talking about randomness, note that the two most important parameters of a random variable are the expectation, or mean value, and the variance. For the function  $F(\sqrt{2}; \beta; \gamma; N)$ , the estimate (4.31) gives the expectation.

*Explaining why the natural scaling is exponential.* Note that for any  $1 < M < N$ , the counting function is slowly changing in the sense that

$$F(\sqrt{2}; \beta; \gamma; N) - F(\sqrt{2}; \beta; \gamma; M) = O(\log(N/M)), \tag{4.35}$$

where  $c_4 \log(N/M)$  is the corresponding area. The geometric reason behind this is the exponentially sparse occurrence of lattice points in the corresponding long and narrow tilted hyperbola. The proof of (4.35) is a straightforward application of Lemma 4 below.

We have the following corollary of (4.35). If  $M = cN$ , i.e.  $n$  runs through the interval  $cN < n < N$  with some constant  $0 < c < 1$ , then the fluctuation of  $F(\sqrt{2}; \beta; \gamma; N)$  is a trivial  $O(1)$ . This negligible constant size change  $O(1)$  in (4.35), as  $n$  runs through  $cN < n < N$ , explains why it is more natural to switch to the exponential scaling  $F(\sqrt{2}; \beta; \gamma; e^N)$ . In the rest of this discussion, we shall often prefer the exponential scaling.

The variance comes from the following non-trivial result. For any  $\gamma > 0$ , there is a positive effective constant  $\sigma = \sigma(\gamma) > 0$  such that

$$\lim_{N \rightarrow \infty} \frac{1}{N} \int_0^1 \left( F(\sqrt{2}; \beta; \gamma; e^N) - \frac{\gamma}{\sqrt{2}} N \right)^2 d\beta = \sigma^2(\gamma).$$

The proof of this limit formula is based on a combination of Fourier analysis (Poisson summation formula, Parseval formula) and the arithmetic of the quadratic number field  $\mathbf{Q}(\sqrt{2})$ ; see [2].

The first probabilistic result, nicely fitting the general scheme of *determinism vs. randomness*, is the following; for the proof, see [2].

**Theorem A (Central Limit Theorem).** *The renormalized counting function*

$$\frac{F(\sqrt{2}; \beta; \gamma; e^N) - (\gamma/\sqrt{2})N}{\sigma(\gamma)\sqrt{N}}, \quad 0 \leq \beta < 1,$$

has a standard normal limit distribution as  $N \rightarrow \infty$ .

To give at least some vague intuition behind Theorem A, we write

$$G_j(\beta) = F(\sqrt{2}; \beta; \gamma; e^j) - F(\sqrt{2}; \beta; \gamma; e^{j-1}), \quad j = 1, 2, \dots, N.$$

In other words,  $G_j(\beta)$  is the number of integral solutions  $n \in \mathbf{N}$  of (4.29) satisfying  $e^{j-1} < n \leq e^j$ .

Note that  $G_j(\beta)$  is a bounded function. This follows from Lemma 4 below, and from the obvious geometric fact that any short hyperbola segment corresponding to  $G_j$  is basically a rectangle. More precisely, any short hyperbola segment corresponding to  $G_j$  can be approximated by an inscribed rectangle  $R_1$  of slope  $1/\sqrt{2}$  and a circumscribed rectangle  $R_2$  of slope  $1/\sqrt{2}$  such that the ratio of the two areas is uniformly bounded by an absolute constant.

It is time now to formulate

**Lemma 4.** *Every tilted rectangle of slope  $1/\sqrt{2}$  and area  $\frac{1}{5}$  contains at most one lattice point.*

We postpone the proof of this simple but important result to the next section.

Lemma 4 can be easily generalized. The same proof gives that for any quadratic irrational  $\alpha$ , there is a positive constant  $c_5 = c_5(\alpha) > 0$  such that every tilted rectangle of slope  $\alpha$  and area  $c_5$  contains at most one lattice point.

Our key intuition is that the bounded function  $G_j(\beta)$  resembles the  $j$ -th Rademacher function, so the sum

$$F(\sqrt{2}; \beta; \gamma; e^N) - \frac{\gamma}{\sqrt{2}}N = \sum_{j=1}^N \left( G_j(\beta) - \frac{\gamma}{\sqrt{2}} \right),$$

as a function of  $\beta \in [0, 1)$ , behaves like a sum of  $N$  independent Bernoulli variables

$$F(\sqrt{2}; \beta; \gamma; e^N) - \frac{\gamma}{\sqrt{2}}N \approx \underbrace{\pm 1 \pm 1 \pm \dots \pm 1}_N,$$

referred to often as an  $N$ -step random walk.

Our next result, Theorem B, can be interpreted as a variant of Khintchine’s famous Law of the Iterated Logarithm in probability theory; see [21]. We show that the number of solutions  $F(\sqrt{2}; \beta; \gamma; e^n)$  of (4.29) oscillates between the sharp bounds

$$\begin{aligned} \frac{\gamma}{\sqrt{2}}n - \sigma \sqrt{n} \sqrt{(2 + \varepsilon) \log \log n} &< F(\sqrt{2}; \beta; \gamma; e^n) \\ &< \frac{\gamma}{\sqrt{2}}n + \sigma \sqrt{n} \sqrt{(2 + \varepsilon) \log \log n}, \end{aligned} \quad (4.36)$$

where  $\varepsilon > 0$ , as  $n \rightarrow \infty$  for almost all  $\beta$ . Note that (4.36) fails with  $2 - \varepsilon$  in place of  $2 + \varepsilon$ , where  $\varepsilon > 0$ . Here the main term  $(\gamma/\sqrt{2})n$  means the area, so (4.36) can be considered a highly sophisticated justification of the Naive Area Principle.

The estimate (4.36) is particularly interesting in view of the fact that the classical Circle Problem is unsolved, and seems to be hopeless by current techniques. What (4.36) means is that we can solve a *Hyperbola Problem* instead of the Circle Problem. More precisely, we can prove for long and narrow tilted hyperbola segments what nobody can prove for large concentric circles. Namely, we can show that for almost all centers, i.e. for almost all values of the translation parameter  $\beta$ , the number of lattice points asymptotically equals the area plus an error which, even in the worst case scenario, is about the square root of the area. For circles the corresponding maximum error should be the square root of the circumference.

The Law of the Iterated Logarithm is one of the most famous results in classical probability theory, and describes the maximum fluctuation in the infinite one-dimensional random walk. The term *infinite random walk* refers to an infinite sequence of random Bernoulli trials, where each trial is tossing a fair coin. Of course, coin tossing belongs to the physical world; it is not a mathematical concept. But there is a well-known pure mathematical problem, which is

considered equivalent. We can study the digit distribution of a typical real number written in binary form

$$\beta = \frac{b_1}{2} + \frac{b_2}{2^2} + \frac{b_3}{2^3} + \dots,$$

where each  $b_i = 0$  or  $1$ ; here we have assumed for simplicity that  $0 < \beta < 1$ . The infinite 0-1 sequence

$$b_1 = b_1(\beta), b_2 = b_2(\beta), b_3 = b_3(\beta), \dots,$$

i.e. the sequence of binary digits of  $0 < \beta < 1$ , represents an infinite heads-and-tails sequence, say, with 1 as heads and 0 as tails. The sum

$$B_n = B_n(\beta) = b_1 + b_2 + b_3 + \dots + b_n$$

counts the number of 1's, or heads, among the first  $n$  binary digits of  $0 < \beta < 1$ . Borel's classical theorem about normal numbers asserts that

$$\frac{B_n(\beta)}{n} \rightarrow \frac{1}{2} \quad \text{for almost all } 0 < \beta < 1.$$

Let  $S_n = S_n(\beta)$  denote the corresponding error term

$$S_n = S_n(\beta) = 2B_n(\beta) - n = \text{number of heads} - \text{number of tails},$$

so that  $S_n = S_n(\beta)$  represents the number of heads minus the number of tails among the first  $n$  random trials, or coin tosses.

A well-known theorem of Khintchine [21] asserts that

$$\limsup_{n \rightarrow \infty} \frac{S_n(\beta)}{\sqrt{2n \log \log n}} = 1 \quad \text{for almost all } 0 < \beta < 1.$$

Note that Khintchine's Theorem is a far-reaching quantitative improvement on Borel's famous theorem on normal numbers. The long form of Khintchine's Theorem says that for any  $\varepsilon > 0$  and almost all  $\beta$ , we have the following two statements:

- $S_n(\beta) < (1 + \varepsilon)\sqrt{2n \log \log n}$  for all sufficiently large values of  $n$ ; and
- $S_n(\beta) > (1 - \varepsilon)\sqrt{2n \log \log n}$  for infinitely many values of  $n$ .

This strikingly elegant and precise result is the simplest form of the so-called Law of the Iterated Logarithm, usually called Khintchine's form.

Let us return to (4.36). The fact that it is an analog of Khintchine's Law of the Iterated Logarithm suggests the vague intuition that the lattice point counting



function  $F(\sqrt{2}; \beta; \gamma; e^n)$  behaves like a generalized digit sum as  $\beta$  runs through  $0 < \beta < 1$ .

What we are going to actually formulate below are two generalizations or refinements of (4.36); see Theorem B. The first generalization is that for almost all  $\beta$ , (4.36) holds for all  $\gamma$ , or in general, for all intervals  $[\gamma_1, \gamma_2]$ . This is a variant of the so-called Cassels’s form of the Law of the Iterated Logarithm; see [12].

The second generalization of (4.36) is the Kolmogorov–Erdős form, an ultimate convergence-divergence criterion, which contains Khintchine’s form as a simple corollary; see [14, 15, 22].

**Theorem B (Law of the Iterated Logarithm).**

(i) Let  $\varepsilon > 0$  be an arbitrarily small but fixed constant. Then for almost all  $\beta$ ,

$$\begin{aligned} \frac{\gamma}{\sqrt{2}}n - \sigma \sqrt{(2 + \varepsilon)n \log \log n} &< F(\sqrt{2}; \beta; \gamma; e^n) \\ &< \frac{\gamma}{\sqrt{2}}n + \sigma \sqrt{(2 + \varepsilon)n \log \log n} \end{aligned} \quad (4.37)$$

holds for all  $\gamma > 0$  and for all sufficiently large  $n$ , i.e. for all  $n > n_0(\beta, \gamma)$ .

(ii) Let  $\varphi(n)$  be an arbitrary positive increasing function of  $n$ . Let  $\gamma > 0$  be fixed. Then for almost all  $\beta$ ,

$$F(\sqrt{2}; \beta; \gamma; e^n) > \frac{\gamma}{\sqrt{2}}n + \varphi(n)\sigma\sqrt{n}$$

holds for infinitely many values of  $n$  if and only if the series

$$\sum_{n=1}^{\infty} \frac{\varphi(n)}{n} e^{-\varphi^2(n)/2} \quad (4.38)$$

diverges. The same conclusion holds for the other inequality

$$F(\sqrt{2}; \beta; \gamma; e^n) < \frac{\gamma}{\sqrt{2}}n - \varphi(n)\sigma\sqrt{n}.$$

Note that (4.37) is sharp in the sense that  $2 + \varepsilon$  cannot be replaced by  $2 - \varepsilon$ .

*Remarks.* (i) By Lemma 1, we have  $f(\sqrt{2}; \beta; c; N) = F(\sqrt{2}; \beta; \gamma; N) + O(1)$  as  $N \rightarrow \infty$ , where  $c = \gamma/2\sqrt{2}$ . So Lemma 1 implies that Theorems A and B remain true if  $F(\sqrt{2}; \beta; \gamma; N)$  is replaced by the number of solutions  $f(\sqrt{2}; \beta; c; N)$  of the inhomogeneous diophantine inequality (4.28).

(ii) In Theorem B(i), there is a dramatic difference between rational  $\beta$  and almost all  $\beta$ . For every rational  $\beta$ , the counting function has the form

$$F(\sqrt{2}; \beta; \gamma; N) = c(\gamma) \log N + O(1) \quad \text{as } N \rightarrow \infty$$

for all  $\gamma > 0$ , and it remains valid if  $\sqrt{2}$  is replaced by any quadratic irrational. This bounded size fluctuation around the main term  $c \log N$ , which is typically not the area, jumps up considerably. By (4.37), we have square root size fluctuations around the main term, which is the area, so the fluctuations have size the square root of the area, and this holds for almost all  $\beta$  and all  $\gamma > 0$ .

Let us return to (4.36). It is a special case of Theorem B(ii) with

$$\varphi(n) = ((2 \pm \varepsilon) \log \log n)^{1/2}.$$

Indeed, the series (4.38) is divergent or convergent depending on whether we have  $2 + \varepsilon$  or  $2 - \varepsilon$  in the definition of  $\varphi(n)$ .

We can obtain a much more delicate result by choosing a large integer  $k \geq 4$  and writing

$$\varphi(n) = (2 \log_2 n + 2 \log_3 n + 2 \log_4 n + \dots + 2 \log_{k-1} n + (2 \pm \varepsilon) \log_k n)^{1/2}.$$

Beware that here, and here only, we use the space-saving notation  $\log_2 n = \log \log n$ , i.e. it means the iterated logarithm instead of the usual meaning as base 2 logarithm, and in general,  $\log_k n = \log(\log_{k-1} n)$  denotes the  $k$ -times iterated logarithm of  $n$ . With this choice of  $\varphi(n)$ , we have

$$\sum_{n=1}^{\infty} \frac{\varphi(n)}{n} e^{-\varphi^2(n)/2} \approx \sum_n \frac{1}{n \log n \log_2 n \log_3 n \dots \log_{k-1} n (\log_k n)^{1 \pm \varepsilon/2}},$$

which is divergent or convergent depending on whether we have  $2 + \varepsilon$  or  $2 - \varepsilon$  in the definition of  $\varphi(n)$ .

This example clearly illustrates the remarkable precision of Theorem B(ii).

Next we focus on a simple consequence of Theorem B. Let  $c > 0$  be arbitrarily small but fixed. Then by Theorem B, the inhomogeneous diophantine inequality

$$\|n\sqrt{2} - \beta\| < \frac{c}{n} \tag{4.39}$$

has infinitely many integer solutions  $n \geq 1$  for almost all  $\beta$ , in the sense of the Lebesgue measure.

Inequality (4.39) corresponds to the hyperbola segment

$$|y - \beta| < \frac{c}{x}, \quad x \geq 1,$$

where  $\beta$  is fixed, and this has infinite area. But we may go further, and consider smaller regions

$$|y - \beta| < \frac{1}{x \log x}, \quad |y - \beta| < \frac{1}{x \log x \log \log x},$$

and the like. They all have infinite area, since

$$\int_{\epsilon}^N \frac{dx}{x \log x} = \log \log N \quad \text{and} \quad \int_{e^e}^N \frac{dx}{x \log x \log \log x} = \log \log \log N,$$

and the rest all tend to infinity as  $N \rightarrow \infty$ . It is very natural, therefore, to ask the following question.

*Question.* Consider the inequalities

$$\|n\sqrt{2} - \beta\| < \frac{c}{n \log n}, \quad n \geq n_1, \tag{4.40}$$

$$\|n\sqrt{2} - \beta\| < \frac{c}{n \log n \log \log n}, \quad n \geq n_2, \tag{4.41}$$

and so on, where  $0 \leq \beta < 1$  is a fixed constant. Is it true that for almost all  $\beta$ , in the sense of the Lebesgue measure, the inequalities (4.40), (4.41) and the like have infinitely many positive integer solutions  $n$ ?

Well, the answer is affirmative.

**Theorem C (Area Principle for  $\sqrt{2}$ ).** *Let  $\psi(x)$  be any positive decreasing function of the real variable  $x$  satisfying*

$$\sum_{n=1}^{\infty} \psi(n) = \infty. \tag{4.42}$$

*Then the inhomogeneous inequality*

$$\|n\sqrt{2} - \beta\| < \psi(n)$$

*has infinitely many integral solutions for almost all  $0 \leq \beta < 1$ , in the sense of Lebesgue measure.*

Furthermore, there is an interesting generalization of Theorem C where  $\sqrt{2}$  is replaced by any real  $\alpha$ .

To explain this generalization, Theorem D below, we recall the basic question of diophantine approximation. We want to decide whether an inequality

$$\left| \alpha - \frac{p}{q} \right| < \frac{1}{q^2}, \quad \text{or equivalently,} \quad |q\alpha - p| < \frac{1}{q},$$

with integers  $p$  and  $q$ , or more generally, an inequality

$$\|q\alpha\| < \psi(q), \tag{4.43}$$

where  $\psi(q)$  is a positive decreasing function of  $q$ , has infinitely many integral solutions in  $q$ , and if this is the case, to determine the solutions, or at least the asymptotic number of integral solutions.

It is perfectly natural to study the inhomogeneous analog of (4.43), the inequality

$$\|q\alpha - \beta\| < \psi(q), \tag{4.44}$$

where  $\beta$  is an arbitrary fixed real number. Of course, we may assume that  $0 \leq \beta < 1$ .

Is there any connection between the solvability of the homogeneous inequality (4.43) and the inhomogeneous inequality (4.44)? Theorem C is about the special case  $\alpha = \sqrt{2}$ , and it justifies the Naive Area Principle. Recall that the Naive Area Principle is a vague intuition claiming that a nice region of infinite area must contain infinitely many lattice points. We know that the Naive Area Principle is false for the hyperbolic region  $-\frac{1}{2} \leq x^2 - 2y^2 \leq \frac{1}{2}$ , which has infinite area and contains only one lattice point, namely the origin. This Pell inequality is basically equivalent to the diophantine inequality

$$\|q\sqrt{2}\| < \frac{c}{q}, \tag{4.45}$$

with  $c \leq 2^{-5/2}$ , and (4.45) does not have infinitely many integral solutions in  $q$  if the constant  $c < 2^{-5/2}$ .

The failure of the Naive Area Principle for (4.45) is compensated by the success of the Naive Area Principle for the inhomogeneous inequality

$$\|q\sqrt{2} - \beta\| < \psi(q),$$

which has infinitely many integral solution  $q$  for almost all  $\beta$ , provided that  $\psi(x)$  is any positive decreasing function of the real variable  $x$  satisfying (4.42). This is the statement of Theorem C. The next result generalizes the special case  $\alpha = \sqrt{2}$  to arbitrary real  $\alpha$ .

**Theorem D (General Area Principle).** *Let  $\psi(x)$  be any positive decreasing function of the real variable  $x$  satisfying (4.42). For any real number  $\alpha$ , at least one of the following two cases always holds:*

- (i) *The homogeneous inequality (4.43) has infinitely many integral solutions.*
- (ii) *The inhomogeneous inequality (4.44) has infinitely many integral solutions for almost all  $0 \leq \beta < 1$ , in the sense of Lebesgue measure.*

*Remark.* Note that divergence condition (4.42) is necessary. Indeed, if

$$\sum_{n=1}^{\infty} \psi(n) < \infty, \tag{4.46}$$

then the set of pairs  $(\alpha, \beta)$ , for which the inequality (4.44) has infinitely many integral solutions  $q$ , has two-dimensional Lebesgue measure zero. This statement immediately follows from the statement that for every fixed  $\beta$ , the set of  $\alpha$  which satisfy (4.44) for infinitely many  $q$  has Lebesgue measure zero. The second statement has an easy proof as follows. Every such  $\alpha$  in  $0 < \alpha < 1$  is contained in infinitely many intervals of the form

$$\left[ \frac{p + \beta}{q} - \frac{\psi(q)}{q}, \frac{p + \beta}{q} + \frac{\psi(q)}{q} \right]$$

with integers  $q \geq N$  and  $1 \leq p \leq q$ , and the total length of these intervals is less than

$$2 \sum_{q \geq N} \psi(q),$$

which by (4.46) tends to zero as  $N \rightarrow \infty$ . This means that Theorem D is a precise convergence-divergence type result, or we may call it a zero-one law, to borrow a well-known concept from probability theory.

Let us return to the inhomogeneous inequality (4.44). If  $\alpha$  is rational and  $\beta$  is irrational, then (4.44) has only finitely many integral solutions for any  $\psi(q) \rightarrow 0$  as  $q \rightarrow \infty$ . Well, this is trivial. It is less trivial to find an irrational  $\alpha$  and a decreasing function  $\psi(x)$  satisfying (4.42) such that for almost all  $\beta$ , (4.44) has only finitely many integral solutions. We can take any irrational  $0 < \alpha < 1$  with sufficiently large partial quotients in the sense that

$$\alpha = \frac{1}{a_1 + \frac{1}{a_2 + \dots}} = [a_1, a_2, a_3, \dots],$$

where

$$a_k \approx k^{(\log k)^2}, \tag{4.47}$$

and take

$$\psi(q) = \frac{1}{q \log q}. \tag{4.48}$$

Then the denominator  $q_k$  of the  $k$ -th convergent of  $\alpha$  is roughly

$$q_k \approx a_1 a_2 \dots a_k \approx k^{k(\log k)^2}, \tag{4.49}$$

and so

$$\sum_k \frac{1}{\log q_k} = o\left(\sum_k \frac{1}{k(\log k)^3}\right) < \infty.$$

We recall the well-known fact

$$\left| \alpha - \frac{p_k}{q_k} \right| < \frac{1}{q_k q_{k+1}}$$

which implies

$$\left| n\alpha - \frac{np_k}{q_k} \right| < \frac{n}{q_k q_{k+1}}. \quad (4.50)$$

If  $q_k \leq n < q_{k+1}k^{-2}$  and

$$\|n\alpha - \beta\| < \frac{1}{n \log n},$$

then by (4.49) and (4.50), we have

$$\left\| \beta - \frac{np_k}{q_k} \right\| < \frac{1}{k^2 q_k} + \frac{1}{n \log n} < \frac{2}{k(\log k)^3 q_k}. \quad (4.51)$$

If  $q_{k+1}k^{-2} \leq n < q_{k+1}$ , then define the set

$$A_k = \bigcup_n \left[ n\alpha - \frac{1}{n \log n}, n\alpha + \frac{1}{n \log n} \right] \pmod{1}, \quad (4.52)$$

where the summation in (4.52) is extended over all  $n$  with  $q_{k+1}k^{-2} \leq n < q_{k+1}$ . Motivated by (4.51), define the set

$$B_k = \bigcup_{0 \leq j < q_k} \left[ \frac{j}{q_k} - \frac{2}{k(\log k)^3 q_k}, \frac{j}{q_k} + \frac{2}{k(\log k)^3 q_k} \right] \pmod{1}. \quad (4.53)$$

Clearly

$$\sum_k \text{meas}(B_k) \leq \sum_k \frac{4}{k(\log k)^3} < \infty, \quad (4.54)$$

where  $\text{meas}$  denotes the usual Lebesgue measure, and

$$\sum_k \text{meas}(A_k) = O\left(\sum_k \frac{\log(k^2)}{k(\log k)^3}\right) = O\left(\sum_k \frac{1}{k(\log k)^2}\right) < \infty. \quad (4.55)$$

It follows from (4.54) and (4.55) that almost all  $\beta$  are contained in only a finite number of  $A_k$  and in a finite number of  $B_k$ . In view of (4.51)–(4.53), this implies

that for almost all  $\beta$ , the inequality (4.44) has only finitely many integral solutions, where  $\alpha$  and  $\psi$  are defined by (4.47) and (4.48).

For the proofs of Theorems A and B, we refer the reader to the forthcoming book [2]. For the proofs of Theorems C and D, see the recent paper [8]. This section was a detour, or rather a counterpart; the rest of the chapter is about extra large fluctuations, i.e. sophisticated violations of the Naive Area Principle.

The next section is technical, and contains the proofs of Theorem 3 and Lemmas 1–4. The truly interesting new results come later, starting in Sect. 4.4.

### 4.3 Proving Theorem 3 and the Lemmas

*Proof of Lemma 2.* First we establish the estimate (4.31). Consider the hyperbolic needle  $H_N(\gamma) = H_N(\sqrt{2}; \gamma)$ , defined by

$$H_N(\gamma) = \{(x, y) \in \mathbf{R}^2 : -\gamma \leq x^2 - 2y^2 \leq \gamma, 1 \leq x + y\sqrt{2} \leq 2\sqrt{2}N\}. \quad (4.56)$$

Comparing (4.11) with (4.56), we see that

$$H_N(\gamma) = H^*(\sqrt{2}; [-\gamma, \gamma]; N),$$

so by (4.13), we deduce that

$$\text{area}(H_N(\gamma)) = \frac{\gamma}{\sqrt{2}} \log N + O(1). \quad (4.57)$$

Next we need the following almost trivial result.

**Lemma 5.** *Let  $A \subset \mathbf{R}^2$  be a Lebesgue measurable set in the plane with finite measure denoted by  $\text{area}(A)$ . Then*

$$\int_0^1 \int_0^1 |(A + \mathbf{x}) \cap \mathbf{Z}^2| \, d\mathbf{x} = \text{area}(A),$$

where  $A + \mathbf{x}$  denotes the translation of the set  $A$  by the vector  $\mathbf{x} \in \mathbf{R}^2$ .

Now by Lemma 5, we have

$$\int_0^1 \int_0^1 |(H_N(\gamma) + \mathbf{v}) \cap \mathbf{Z}^2| \, d\mathbf{v} = \text{area}(H_N(\gamma)). \quad (4.58)$$

If  $\mathbf{v} = (v_1, v_2) \in [0, 1)^2$  is chosen in such a way that  $v_1 - v_2\sqrt{2} \equiv \beta \pmod 1$  is fixed, then clearly

$$|F(\sqrt{2}; \beta; \gamma; N) - |(H_N(\gamma) + \mathbf{v}) \cap \mathbf{Z}^2|| < c_6(\gamma), \quad (4.59)$$

where  $c_6(\gamma)$  is a constant independent of  $\beta$  and  $N$ . The estimate (4.31) follows on combining (4.57)–(4.59).

Next we prove (4.32). Let  $0 \leq a < b \leq 1$  be fixed. For any  $M \geq 1$ , consider the parallelogram

$$\mathcal{P}_M = \{\mathbf{v} = (v_1, v_2) \in \mathbf{R}^2 : a \leq v_1 - v_2\sqrt{2} \leq b, 0 \leq v_1 + v_2\sqrt{2} \leq M\}.$$

If  $M$  is large, then  $\mathcal{P}_M$  is a long and narrow parallelogram, but we can then turn it into a *round* shape by applying an appropriate automorphism of the quadratic form  $x^2 - 2y^2$ . The substitution  $x_1 = x + 2y$ ,  $y_1 = x + y$  is a fundamental automorphism,<sup>5</sup> and writing

$$A = \begin{pmatrix} 1 & 2 \\ 1 & 1 \end{pmatrix},$$

we note that  $A^k$ ,  $k \in \mathbf{Z}$ , give rise to infinitely many automorphisms preserving the lattice points and the area. The eigenvectors of the matrix  $A$  are parallel to the sides of parallelogram  $\mathcal{P}_M$ , so on applying an appropriate power  $A^k$  on the long and narrow parallelogram  $\mathcal{P}_M$ , we obtain a *round* parallelogram  $A^k \mathcal{P}_M$  with sides parallel to that of  $\mathcal{P}_M$ , and

$$\text{area}(A^k \mathcal{P}_M) = \text{area}(\mathcal{P}_M) = c_7 M.$$

Here *round* means that the diameter of parallelogram  $A^k \mathcal{P}_M$  is  $O(\sqrt{M})$ , so the number of unit squares  $[0, 1)^2 + \mathbf{n}$ ,  $\mathbf{n} \in \mathbf{Z}^2$ , intersecting the boundary of  $A^k \mathcal{P}_M$  is  $O(\sqrt{M})$ .

Combining this geometric fact with (4.58), we have

$$\frac{1}{\text{area}(\mathcal{P}_M)} \int_{\mathcal{P}_M} |(H_N(\gamma) + \mathbf{v}) \cap \mathbf{Z}^2| \, d\mathbf{v} = \text{area}(H_N(\gamma))(1 + O(M^{-1/2})). \tag{4.60}$$

If  $\mathbf{v} = (v_1, v_2) \in [0, 1)^2$  is chosen in such a way that  $v_1 - v_2\sqrt{2} \equiv \beta \pmod{1}$  is fixed, then clearly

$$|F(\sqrt{2}; \beta; \gamma; N) - |(H_N(\gamma) + \mathbf{v}) \cap \mathbf{Z}^2|| < c_8(\gamma, M), \tag{4.61}$$

where  $c_8(\gamma, M)$  is a constant independent of  $\beta$  and  $N$ . Combining (4.57), (4.60) and (4.61), we have

$$\frac{\frac{1}{b-a} \int_a^b F(\sqrt{2}; \beta; \gamma; N) \, d\beta}{\log N}$$

---

<sup>5</sup>Indeed, we have  $x_1^2 - 2y_1^2 = (x + 2y)^2 - 2(x + y)^2 = -(x^2 - 2y^2)$ .



$$= \left( \frac{\gamma}{\sqrt{2}} + O\left(\frac{1}{\log N}\right) \right) (1 + O(M^{-1/2})) + \frac{c_8(\gamma, M)}{\log N}. \quad (4.62)$$

Since  $M$  can be arbitrarily large, (4.62) implies (4.32). The proof of Lemma 2 is now complete.  $\square$

*Proof of Lemma 5.* First assume that  $A$  is bounded. Let  $N$  be a large integer. In view of the periodicity of  $\mathbf{Z}^2$ , we have

$$\int_0^N \int_0^N |(A + \mathbf{x}) \cap \mathbf{Z}^2| \, d\mathbf{x} = N^2 \int_0^1 \int_0^1 |(A + \mathbf{x}) \cap \mathbf{Z}^2| \, d\mathbf{x}.$$

On the other hand,

$$\begin{aligned} \int_0^N \int_0^N |(A + \mathbf{x}) \cap \mathbf{Z}^2| \, d\mathbf{x} &= \sum_{\mathbf{n} \in \mathbf{Z}^2} \text{area}\{\mathbf{x} \in [0, N]^2 : \mathbf{n} \in A + \mathbf{x}\} \\ &= \sum_{\mathbf{n} \in \mathbf{Z}^2} \text{area}\{(\mathbf{n} - A) \cap [0, N]^2\}. \end{aligned}$$

Without loss of generality, we can assume that the origin is inside  $A$ . Let  $d(A)$  denote the diameter of  $A$ . Then  $(\mathbf{n} - A) \subset [0, N]^2$  if  $\mathbf{n} \in [d(A), N - d(A)]^2$ . On the other hand,  $(\mathbf{n} - A) \cap [0, N]^2 = \emptyset$  if  $\mathbf{n} \notin [-d(A), N + d(A)]^2$ . Thus we have

$$(N + 2d(A))^2 \cdot \text{area}(A) \geq \sum_{\mathbf{n} \in \mathbf{Z}^2} \text{area}\{(\mathbf{n} - A) \cap [0, N]^2\} \geq (N - 2d(A))^2 \cdot \text{area}(A).$$

Dividing the last inequalities by  $N^2$ , and combining with the equations above, we see that Lemma 5 follows as  $N$  tends to infinity. If  $A$  is unbounded, then we approximate  $A$  by an increasing sequence  $A_1 \subset A_2 \subset A_3 \subset \dots$  of subsets of  $A$  such that each  $A_k$  is bounded and  $\text{area}(A \setminus A_k) \rightarrow 0$ . The last step is then to use the continuity of the Lebesgue measure.  $\square$

*Proof of Lemma 1.* For notational simplicity, we restrict our proof to the special case  $\beta_2 = 0$ ; the general case is the same. Again the key step is to apply Lemma 5. For  $1 \leq K < L \leq \infty$ , consider the four regions

$$\begin{aligned} H_{K,L}(\beta; \gamma) &= \{(x, y) \in \mathbf{R}^2 : -\gamma \leq (x + \beta)^2 - 2y^2 \leq \gamma, K \leq y \leq L, x > 0\}, \\ \tilde{H}_{K,L}(\beta; \gamma) &= \{(x, y) \in \mathbf{R}^2 : 2\sqrt{2}y|x + \beta - y\sqrt{2}| < \gamma, K \leq y \leq L, x > 0\}, \\ \tilde{H}_{K,L}^+(\beta; \gamma) &= \{(x, y) \in \mathbf{R}^2 : (2\sqrt{2}y + 1)|x + \beta - y\sqrt{2}| < \gamma, K \leq y \leq L, x > 0\}, \\ \tilde{H}_{K,L}^-(\beta; \gamma) &= \{(x, y) \in \mathbf{R}^2 : (2\sqrt{2}y - 1)|x + \beta - y\sqrt{2}| < \gamma, K \leq y \leq L, x > 0\}. \end{aligned}$$

In view of the factorization (4.27), the condition  $(x, y) \in H_{K,L}(\beta; \gamma)$  gives the estimate  $x + \beta = y\sqrt{2} + o(1)$ . In fact, we have the stronger form  $x + \beta = y\sqrt{2} + O(1/y)$ . Thus there is a threshold  $c_9 = c_9(\gamma)$  such that

$$\tilde{H}_{K,L}^+(\beta; \gamma) \subset H_{K,L}(\beta; \gamma) \subset \tilde{H}_{K,L}^-(\beta; \gamma)$$

for all  $L > K > c_9(\gamma)$ . On the other hand, it is trivial that

$$\tilde{H}_{K,L}^+(\beta; \gamma) \subset \tilde{H}_{K,L}(\beta; \gamma) \subset \tilde{H}_{K,L}^-(\beta; \gamma).$$

Consider now the special case  $K = 1, L = \infty, \beta = 0$ , and study the difference set

$$D(\gamma) = \tilde{H}_{1,\infty}^-(0; \gamma) \setminus \tilde{H}_{1,\infty}^+(0; \gamma).$$

The area of this difference set can be estimated by

$$\begin{aligned} \text{area}(D(\gamma)) &= O\left(\int_1^\infty \left(\frac{1}{2\sqrt{2}y-1} - \frac{1}{2\sqrt{2}y+1}\right) dy\right) \\ &= O\left(\int_1^\infty \frac{dy}{8y^2-1}\right) = O(1). \end{aligned}$$

Combining this with Lemma 5, we have

$$\int_0^1 \int_0^1 |(D(\gamma) + \mathbf{v}) \cap \mathbf{Z}^2| d\mathbf{v} = \text{area}(D(\gamma)) < \infty. \tag{4.63}$$

If  $\mathbf{v} = (v_1, v_2) \in [0, 1)^2$  is chosen in such a way that  $v_1 - v_2\sqrt{2} \equiv \beta \pmod 1$  is fixed, then

$$D(\gamma) + \mathbf{v} \supset H_{K,L}(\beta; \gamma) \Delta \tilde{H}_{K,L}^+(\beta; \gamma), \tag{4.64}$$

where  $A \Delta B = (A \setminus B) \cup (B \setminus A)$  denotes the symmetric difference of the sets  $A$  and  $B$ . Combining (4.63) and (4.64), Lemma 1 follows easily.  $\square$

*Proof of Lemma 4.* Consider a rectangle of slope  $1/\sqrt{2}$  which contains two lattice points  $P = (k, \ell)$  and  $Q = (m, n)$ ; in fact, assume that  $P, Q$  are two vertices of the rectangle. We denote the vector from  $P$  to  $Q$  by  $\mathbf{v} = (m - k, n - \ell)$ , and consider the two perpendicular unit vectors

$$\mathbf{e}_1 = \left(\frac{\sqrt{2}}{\sqrt{3}}, \frac{1}{\sqrt{3}}\right) \quad \text{and} \quad \mathbf{e}_2 = \left(\frac{1}{\sqrt{3}}, -\frac{\sqrt{2}}{\sqrt{3}}\right).$$

Then the two side lengths  $a$  and  $b$  of the rectangle can be expressed in terms of the inner products

$$a = |\mathbf{e}_1 \cdot \mathbf{v}| = \frac{|p\sqrt{2} + q|}{\sqrt{3}} \quad \text{and} \quad b = |\mathbf{e}_2 \cdot \mathbf{v}| = \frac{|p - q\sqrt{2}|}{\sqrt{3}},$$

where  $p = m - k$  and  $q = n - \ell$ . Thus we have

$$\text{area} = ab = \frac{|(p\sqrt{2} + q)(p - q\sqrt{2})|}{3}.$$

Without loss of generality, we can assume that  $p \geq 0$  and  $q \geq 0$ . Since  $(p, q) \neq (0, 0)$ , we have  $|p - q\sqrt{2}| = 1/(p + q\sqrt{2})$ , and so

$$\text{area} = \frac{|(p\sqrt{2} + q)(p - q\sqrt{2})|}{3} = \frac{p\sqrt{2} + q}{3(p + q\sqrt{2})} \geq \frac{p + q}{3(p\sqrt{2} + q\sqrt{2})} = \frac{1}{3\sqrt{2}} > \frac{1}{5},$$

proving Lemma 4. □

*Proof of Theorem 3.* We shall show that the set of numbers  $\beta$  in question, the set of divergence points, contains a Cantor set. This guarantees that the cardinality of the set is continuum.

We make a standard Cantor set construction, i.e. we apply the method of nested intervals. For notational convenience, we write  $F(\sqrt{2}; \beta; \gamma; N) = F(\beta; \gamma; N)$ . By (4.31), we have

$$\int_0^1 F(\beta; \gamma; N) \, d\beta = \frac{\gamma}{\sqrt{2}} \log N + O(1).$$

Applying this with  $\gamma = \frac{1}{4}$ , we obtain the existence of  $0 < \beta_1 < 1$  and an arbitrarily large integer  $N_1$  such that

$$F(\beta_1; \gamma = 1/4; N_1) > \frac{1}{8} \log N_1.$$

Since  $\frac{1}{4} < \frac{1}{2}$ , there exists an interval  $I_1 = [a, b]$  with  $0 < a < b < 1$  such that  $\beta_1 \in I_1$  and

$$F(\beta; \gamma = 1/2; N_1) > \frac{1}{8} \log N_1 \quad \text{for all } \beta \in I_1. \tag{4.65}$$

Next let  $\mathbf{n} = (n_1, n_2) \in \mathbf{Z}^2$  be a lattice point such that  $\beta_2 = n_1 - n_2\sqrt{2} \in I_1$ . Since the equation  $|x^2 - 2y^2| \leq \frac{3}{4}$  does not have a non-zero integral solution, trivially

$$F(\beta_2; \gamma = 3/4; N) < \frac{1}{100} \log N \quad \text{for all } N \geq N_2,$$

where  $N_2$  is a sufficiently large threshold. We can clearly assume that  $N_2 > N_1$ . Since  $\frac{3}{4} > \frac{1}{2}$ , there exists<sup>6</sup> an interval  $I_2 = [a, b]$  with some  $0 < a < b < 1$  such that  $\beta_2 \in I_2$  and

$$F(\beta; \gamma = 1/2; N_2) < \frac{1}{100} \log N_2 \quad \text{for all } \beta \in I_2. \tag{4.66}$$

We can clearly assume that  $I_2$  is a proper subinterval of  $I_1$ . Let  $I(0) = I_2$ . Repeating the second argument, we deduce that there exists another closed subinterval  $I(1)$  such that  $I(0)$  and  $I(1)$  are disjoint,  $I(0) \cup I(1) \subset I_1$  and

$$F(\beta; \gamma = 1/2; N_2^{(1)}) < \frac{1}{100} \log N_2^{(1)} \quad \text{for all } \beta \in I(1). \tag{4.67}$$

We can clearly assume that  $N_2^{(1)} > N_1$ .

By (4.32), we have

$$\frac{1}{|I(0)|} \int_{I(0)} F(\beta; \gamma; N) \, d\beta = (1 + o(1)) \frac{\gamma}{\sqrt{2}} \log N,$$

and applying this with  $\gamma = \frac{1}{4}$ , we obtain the existence of  $0 < \beta_3 < 1$  and a large integer  $N_3$  such that

$$F(\beta_3; \gamma = 1/4; N_3) > \frac{1}{8} \log N_3.$$

Since  $\frac{1}{4} < \frac{1}{2}$ , there exists an interval  $I_3 = [a, b]$  with  $0 < a < b < 1$  such that  $\beta_3 \in I_3$  and

$$F(\beta; \gamma = 1/2; N_3) > \frac{1}{8} \log N_3 \quad \text{for all } \beta \in I_3. \tag{4.68}$$

We can clearly assume that  $I_3$  is a proper subinterval of  $I(0)$ . Write  $I(0, 0) = I_3$ . Similarly, there exists another subinterval  $I(0, 1)$  such that  $I(0, 0)$  and  $I(0, 1)$  are disjoint,  $I(0, 0) \cup I(0, 1) \subset I(0)$  and

$$F(\beta; \gamma = 1/2; N_3^{(1)}) > \frac{1}{8} \log N_3^{(1)} \quad \text{for all } \beta \in I(0, 1). \tag{4.69}$$

There are similar disjoint subintervals  $I(1, 0)$  and  $I(1, 1)$  of  $I(1)$ .

Next, let  $\mathbf{n} = (n_1, n_2) \in \mathbf{Z}^2$  be a lattice point such that  $\beta_4 = n_1 - n_2 \sqrt{2} \in I(0, 0)$ . Since the inequality  $|x^2 - 2y^2| \leq \frac{3}{4}$  does not have a non-trivial integral solution,

$$F(\beta_4; \gamma = 3/4; N) < \frac{1}{100} \log N \quad \text{for all } N \geq N_4,$$

---

<sup>6</sup>Here  $a$  and  $b$  are generic numbers.

where  $N_4 < \infty$  is a sufficiently large threshold. We can clearly assume that  $N_4 > N_3$ . Since  $\frac{3}{4} > \frac{1}{2}$ , there exists an interval  $I_4 = [a, b]$  with  $0 < a < b < 1$  such that  $\beta_4 \in I_4$  and

$$F(\beta; \gamma = 1/2; N_4) < \frac{1}{100} \log N_4 \quad \text{for all } \beta \in I_4. \tag{4.70}$$

We can clearly assume that  $I_4$  is a proper subinterval of  $I(0, 0)$ . Let  $I(0, 0, 0) = I_4$ . Repeating the last argument, there exists another closed subinterval  $I(0, 0, 1)$  such that  $I(0, 0, 0)$  and  $I(0, 0, 1)$  are disjoint,  $I(0, 0, 0) \cup I(0, 0, 1) \subset I(0, 0)$  and

$$F(\beta; \gamma = 1/2; N_4^{(1)}) < \frac{1}{100} \log N_4^{(1)} \quad \text{for all } \beta \in I(0, 0, 1), \tag{4.71}$$

and so on. Repeating this argument, we build an infinite binary tree

$$I_1 \supset I_{\varepsilon_1} \supset I_{\varepsilon_1, \varepsilon_2} \supset I_{\varepsilon_1, \varepsilon_2, \varepsilon_3} \supset \dots,$$

where  $\varepsilon_1, \varepsilon_2, \varepsilon_3, \dots \in \{0, 1\}$ .

For an arbitrary infinite 0-1 sequence  $\varepsilon_1, \varepsilon_2, \varepsilon_3, \dots$ , let

$$\beta \in I_1 \cap I_{\varepsilon_1} \cap I_{\varepsilon_1, \varepsilon_2} \cap I_{\varepsilon_1, \varepsilon_2, \varepsilon_3} \cap \dots$$

Then by (4.65)–(4.71), there exists an infinite sequence  $1 < M_1 < M_2 < M_3 < M_4 < \dots$  of integers such that

$$F(\beta; \gamma = 1/2; M_{2k-1}) > \frac{1}{8} \log M_{2k-1} \quad \text{and} \quad F(\beta; \gamma = 1/2; M_{2k}) < \frac{1}{100} \log M_{2k},$$

where  $k = 1, 2, 3, \dots$ . This proves Theorem 3. □

## 4.4 The Riesz Product and Theorem 12

### 4.4.1 The Method of Nested Intervals vs. the Riesz Product

At the end of Sect. 4.1, we formulated a far-reaching generalization of Theorem 3; see (4.34). It states that Theorem 3 actually holds for every  $\gamma > 0$ , and we have the stronger inequality

$$\limsup_{n \rightarrow \infty} \frac{F(\sqrt{2}; \beta^*; \gamma; n)}{\log n} > \frac{\gamma}{\sqrt{2}} > \liminf_{n \rightarrow \infty} \frac{F(\sqrt{2}; \beta^*; \gamma; n)}{\log n}, \tag{4.72}$$

where  $(\gamma/\sqrt{2}) \log n + O(1)$  is the area of the corresponding hyperbolic region. Indeed, (4.72) holds for continuum many divergence points  $\beta^* = \beta^*(\gamma) \in [0, 1)$ .

The proof of Theorem 3 was based on an elementary argument that we may call the method of nested intervals. To prove (4.72), we need a new idea, and apply a more sophisticated Riesz product argument. The Riesz product is a powerful tool in Fourier analysis. A typical application is to prove large fluctuations for lacunary trigonometric series. To compare the method of nested intervals to the method of Riesz products, we give a simple illustration; see Facts 1 and 2 below.

Consider a finite cosine sum

$$F(x) = \sum_{j=1}^N a_j \cos(2\pi n_j x), \quad \text{where } a_j = \pm 1 \text{ for all } 1 \leq j \leq N, \quad (4.73)$$

and  $1 \leq n_1 < n_2 < \dots < n_N$  are integers. We study the following question. What can we say about  $\max_{0 \leq x \leq 1} F(x)$ ? Well, under different extra conditions, we have different results. We begin with

**Fact 6.** *If the strong gap condition  $n_{j+1}/n_j \geq 8$  holds for every  $1 \leq j \leq N - 1$ , then*

$$\max_{0 \leq x \leq 1} F(x) \geq \frac{N}{2}.$$

*Proof.* The proof is almost trivial. Let

$$J_1 = \left\{ x \in [0, 1] : \cos(2\pi n_1 x) \text{ lies between } \frac{a_1}{2} \text{ and } a_1 \right\}.$$

Since  $a_1 = \pm 1$ , the set  $J_1$  contains a closed subinterval  $I_1$  of length  $|I_1| \geq 1/4n_1$ . Next let

$$J_2 = \left\{ x \in I_1 : \cos(2\pi n_2 x) \text{ lies between } \frac{a_2}{2} \text{ and } a_2 \right\}.$$

Since  $a_2 = \pm 1$ , the set  $J_2$  contains a closed subinterval  $I_2$  of length  $|I_2| \geq 1/4n_2$ . Next let

$$J_3 = \left\{ x \in I_2 : \cos(2\pi n_3 x) \text{ lies between } \frac{a_3}{2} \text{ and } a_3 \right\},$$

and so on. Repeating this process  $N$  times, we obtain a nested sequence of closed intervals

$$[0, 1] \supset I_1 \supset I_2 \supset \dots \supset I_N$$

such that  $a_k \cos(2\pi n_k x) \geq \frac{1}{2}$  for all  $x \in I_k, k = 1, 2, \dots, N$ . Then clearly  $F(x) \geq N/2$  for every  $x \in I_N$ . □

This is a typical application of the method of nested intervals. Next comes the Riesz product argument. The problem that we study is the following. What will happen if the strong gap condition  $n_{j+1}/n_j \geq 8$  is replaced by the weaker condition  $n_{j+1}/n_j \geq 1 + \varepsilon > 1$ , where  $\varepsilon > 0$  is an arbitrarily small but fixed constant? Can we still prove a linear lower bound like  $\max_{0 \leq x \leq 1} F(x) \geq cN$  with some constant  $c = c(\varepsilon) > 0$  depending only on the value of  $\varepsilon$ ? Unfortunately, the method of nested intervals hopelessly collapses. Our new approach is the Riesz product argument. The following result, a well-known theorem of Sidon in Fourier analysis, is much deeper than Fact 6.

**Fact 7 (Sidon’s Theorem).** *If the weak gap condition*

$$\frac{n_{j+1}}{n_j} \geq 1 + \varepsilon > 1 \tag{4.74}$$

*holds for every  $1 \leq j \leq N - 1$ , where  $0 < \varepsilon < \frac{1}{2}$  is a fixed constant, then for  $F(x)$  defined in (4.73), we have*

$$\max_{0 \leq x \leq 1} F(x) \geq cN \quad \text{with} \quad c = \frac{1}{4\varepsilon^{-1} \log(2\varepsilon^{-1})}.$$

*Proof.* Let  $1 = i(1) < i(2) < \dots < i(M)$  be a subsequence of  $1, 2, 3, \dots, N$  such that

$$\frac{n_{i(j+1)}}{n_{i(j)}} \geq \frac{2}{\varepsilon}, \quad j = 1, 2, \dots, M - 1, \tag{4.75}$$

and consider the Riesz product

$$R(x) = \prod_{j=1}^M (1 + a_{i(j)} \cos(2\pi n_{i(j)} x)).$$

Since  $a_{i(j)} = \pm 1$ , we have  $R(x) \geq 0$ . We shall use this Riesz product  $R(x)$  as a test function. First we evaluate the integral

$$\int_0^1 F(x)R(x) dx = \sum_{j=1}^M a_{i(j)}^2 \int_0^1 \cos^2(2\pi n_{i(j)} x) dx = \frac{M}{2}. \tag{4.76}$$

Indeed, multiplying out the Riesz product  $R(x)$ , and then using Euler’s formula  $2e^y = e^{iy} + e^{-iy}$ , we obtain terms like

$$a_{i(j_1)} a_{i(j_2)} a_{i(j_3)} \dots a_{i(j_k)} e^{2\pi i(\pm n_{i(j_1)} \pm n_{i(j_2)} \pm n_{i(j_3)} \pm \dots \pm n_{i(j_k)})}, \tag{4.77}$$

where we shall call (4.77) a product of length  $k \geq 1$ . We distinguish two cases.

Case 8 (short products).  $k = 1$ . Multiplying the corresponding terms with  $F(x)$  and integrating from 0 to 1, we obtain

$$\sum_{j=1}^M a_{i(j)}^2 \int_0^1 \cos^2(2\pi n_{i(j)}x) dx = \frac{M}{2},$$

which is precisely (4.76).

Case 9 (long products).  $k \geq 2$ . We can clearly write  $1 \leq j_1 < j_2 < \dots < j_k$ . Then using the elementary inequalities

$$1 + \frac{\varepsilon}{2} + \left(\frac{\varepsilon}{2}\right)^2 + \left(\frac{\varepsilon}{2}\right)^3 + \dots < 1 + \varepsilon \quad \text{and} \quad 1 - \frac{\varepsilon}{2} - \left(\frac{\varepsilon}{2}\right)^2 - \left(\frac{\varepsilon}{2}\right)^3 - \dots > \frac{1}{1 + \varepsilon}$$

if  $0 < \varepsilon < \frac{1}{2}$ , we deduce that

$$|\pm n_{i(j_1)} \pm n_{i(j_2)} \pm n_{i(j_3)} \pm \dots \pm n_{i(j_k)}| \text{ lies between } (1 + \varepsilon)n_{i(j_k)} \text{ and } \frac{1}{1 + \varepsilon}n_{i(j_k)}.$$

Comparing this to the gap condition (4.74), we see that  $F(x)$  and the long products of  $R(x)$  represent disjoint sets of exponential functions

$$e^{2\pi i \ell x}, \quad \ell \in \mathbf{Z}.$$

Using the orthogonality of these functions, the contribution of Case 9 to the integral  $\int_0^1 F(x)R(x) dx$  is zero. This proves (4.76).

The same argument shows that

$$\int_0^1 R(x) dx = 1. \tag{4.78}$$

Since  $R(x) \geq 0$ , the condition (4.78) means that the integral  $\int_0^1 F(x)R(x) dx$  is a weighted average of  $F(x)$ , with non-negative weights. It follows from (4.76) that

$$\max_{0 \leq x \leq 1} F(x) \geq \int_0^1 F(x)R(x) dx = \frac{M}{2}. \tag{4.79}$$

The inequality  $(1 + \varepsilon)^r > 2/\varepsilon$  clearly holds with  $r = 2\varepsilon^{-1} \log(2\varepsilon^{-1})$ . Thus by (4.74) and (4.75), we can choose

$$M \geq \frac{N}{r} = \frac{N}{2\varepsilon^{-1} \log(2\varepsilon^{-1})}. \tag{4.80}$$

Sidon's theorem then follows from (4.79) and (4.80). □



### 4.4.2 The Rectangle Property and Theorem 12

Let us return now to Theorem 3 and (4.72). We restate Theorem 3 in a slightly different form. Recall the notation in (4.56). We have

$$H_N(\sqrt{2}; \gamma) = \{(x, y) \in \mathbf{R}^2 : -\gamma \leq x^2 - 2y^2 \leq \gamma, 1 \leq x + y\sqrt{2} \leq 2\sqrt{2}N\}, \tag{4.81}$$

that is,  $H_N(\sqrt{2}; \gamma)$  is a long, narrow, tilted hyperbolic needle of slope  $1/\sqrt{2}$ . Its area is  $(\gamma/\sqrt{2}) \log N + O(1)$ ; see (4.57). Theorem 3 states, roughly speaking, that in the special case  $\gamma = \frac{1}{2}$ , there are two translated copies of the same tilted hyperbolic needle  $H_N(\sqrt{2}; \gamma = 1/2)$  such that one is substantially richer in lattice points than the other. The discrepancy is proportional to the area, and we have extra large deviation. More precisely, there is a positive absolute constant  $c_{10} > 0$  such that for infinitely many integers  $N_i$ , where  $N_i \rightarrow \infty$ , there are translated copies  $\mathbf{x}_1^{(i)} + H_{N_i}(\sqrt{2}; \gamma)$  and  $\mathbf{x}_2^{(i)} + H_{N_i}(\sqrt{2}; \gamma)$  of the tilted hyperbolic needle  $H_{N_i}(\sqrt{2}; \gamma = 1/2)$  such that

$$\begin{aligned} & |Z^2 \cap (\mathbf{x}_1^{(i)} + H_{N_i}(\sqrt{2}; \gamma = 1/2))| - |Z^2 \cap (\mathbf{x}_2^{(i)} + H_{N_i}(\sqrt{2}; \gamma = 1/2))| \\ & > c_{10} \log N_i. \end{aligned} \tag{4.82}$$

In view of the periodicity of the lattice points, we can clearly assume that the pairs of vectors  $\mathbf{x}_1^{(i)}$  and  $\mathbf{x}_2^{(i)}$  are all in the unit square  $[0, 1)^2$ , with  $i \rightarrow \infty$ .

The extra large deviation result (4.82), which is equivalent to Theorem 3, can be generalized in several stages. The first generalization is (4.72), or at least an equivalent form as follows.

**Proposition 10.** *Let  $\gamma > 0$  be an arbitrary but fixed real number, and let  $N \geq 2$  be an integer. Then there exists a positive constant  $\delta' = \delta'(\gamma) > 0$ , independent of  $N$ , such that for the tilted hyperbolic needle  $H_N(\sqrt{2}; \gamma)$  of area  $(\gamma/\sqrt{2}) \log N + O(1)$ , there exist translated copies  $\mathbf{x}_1 + H_N(\sqrt{2}; \gamma)$  and  $\mathbf{x}_2 + H_N(\sqrt{2}; \gamma)$  such that*

$$|Z^2 \cap (\mathbf{x}_1 + H_N(\sqrt{2}; \gamma))| > \frac{\gamma}{\sqrt{2}} \log N + \delta' \log N$$

and

$$|Z^2 \cap (\mathbf{x}_2 + H_N(\sqrt{2}; \gamma))| < \frac{\gamma}{\sqrt{2}} \log N - \delta' \log N.$$

Note that Proposition 10 immediately leads to the existence of a single divergence point  $\beta^* = \beta^*(\gamma) \in [0, 1)$  in (4.72). To exhibit continuum many divergence points  $\beta^* = \beta^*(\gamma) \in [0, 1)$ , we simply have to combine Proposition 10 with the routine Cantor set argument in the proof of Theorem 3.

For the second stage of generalization, we replace the set  $\mathbf{Z}^2$  of lattice points in the plane with an arbitrary subset  $\mathcal{A} \subset \mathbf{Z}^2$  of positive density. Here is an illustration of such a set  $\mathcal{A}$ . We say that a lattice point  $\mathbf{n} = (n_1, n_2) \in \mathbf{Z}^2$  is *coprime*<sup>7</sup> if the coordinates  $n_1$  and  $n_2$  are relatively prime. Let  $\mathbf{Z}_{\text{coprime}}^2$  denote the set of coprime lattice points in the plane. It is well known from number theory that  $\mathbf{Z}_{\text{coprime}}^2$  is a subset of  $\mathbf{Z}^2$  with positive density  $6/\pi^2$ .

Now let  $\mathcal{A}$  be an arbitrary subset of  $\mathbf{Z}^2$  of positive density  $\delta = \delta(\mathcal{A}) > 0$ . There is a natural generalization of Proposition 10 where we replace  $\mathbf{Z}^2$  with  $\mathcal{A}$ . The price that we have to pay is that, due to the lack of periodicity of a general subset  $\mathcal{A}$ , the translations are not necessarily in the unit square anymore.

**Proposition 11.** *Let  $\mathcal{A} \subset \mathbf{Z}^2$  be an arbitrary subset of positive density  $\delta = \delta(\mathcal{A}) > 0$ . Let  $\gamma > 0$  be an arbitrary but fixed real number, and let  $N \geq 2$  be an integer. Assume further that  $M/N$  is sufficiently large, depending only on  $\gamma$  and  $\delta$ . Then there exists a positive constant  $\delta' = \delta'(\gamma, \delta) > 0$ , independent of  $N$  and  $M$ , such that for the tilted hyperbolic needle  $H_N(\sqrt{2}; \gamma)$  of area  $(\gamma/\sqrt{2}) \log N + O(1)$ , there exist translated copies  $\mathbf{x}_1 + H_N(\sqrt{2}; \gamma) \subset [0, M]^2$  and  $\mathbf{x}_2 + H_N(\sqrt{2}; \gamma) \subset [0, M]^2$  such that*

$$|\mathcal{A} \cap (\mathbf{x}_1 + H_N(\sqrt{2}; \gamma))| > \frac{\delta\gamma}{\sqrt{2}} \log N + \delta' \log N$$

and

$$|\mathcal{A} \cap (\mathbf{x}_2 + H_N(\sqrt{2}; \gamma))| < \frac{\delta\gamma}{\sqrt{2}} \log N - \delta' \log N.$$

It turns out that the only relevant property of a lattice point set  $\mathcal{A} \subset \mathbf{Z}^2$  that we really use in the proof of Proposition 11 is the rectangle property in Lemma 4, that every tilted rectangle of slope  $1/\sqrt{2}$  and area  $\frac{1}{5}$  contains at most one lattice point. Of course, the concrete value  $\frac{1}{5}$  of the constant is secondary.

The third stage of generalization goes far beyond the family of lattice point sets  $\mathcal{A} \subset \mathbf{Z}^2$ . The only requirement is that the point set satisfies the rectangle property.

**Theorem 12.** *Let  $\mathcal{P}$  be a finite set of points in the square  $[0, M]^2$  with density  $\delta$ , so that the number of elements of  $\mathcal{P}$  is  $|\mathcal{P}| = \delta M^2$ . Assume further that  $\mathcal{P}$  satisfies the following rectangle property, that there is a positive constant  $c_1 = c_1(\mathcal{P}) > 0$  such that every tilted rectangle of slope  $1/\sqrt{2}$  and area  $c_1$  contains at most one element of the set  $\mathcal{P}$ . Let*

---

<sup>7</sup>We also say that such a point is *visible*, explained by the geometric fact that the line segment with  $\mathbf{n}$  and the origin as endpoints does not contain another lattice point. If  $\mathbf{n} = (n_1, n_2) \in \mathbf{Z}^2$  were not coprime, then the point  $(n_1/d, n_2/d) \in \mathbf{Z}^2$ , where  $d \geq 2$  is the greatest common divisor of  $n_1$  and  $n_2$ , would lie on this line segment.

$$\delta' = \delta'(c_1, \gamma, \delta) = 10^{-12}\delta\kappa, \tag{4.83}$$

where

$$\kappa = \min \left\{ \frac{\gamma}{20}, \sqrt{c_1\gamma}, \frac{10^{-7}c_1}{2}, \frac{10^{-7}c_1^2}{2\gamma} \right\}. \tag{4.84}$$

Furthermore, assume that both  $N$  and  $M/N$  are sufficiently large and satisfy

$$N \geq 2^{10(\gamma+\gamma^{-1})} \quad \text{and} \quad M > \frac{10^{11}(\gamma + \gamma^{-1})(N + 2\gamma)}{c_1\delta\kappa}. \tag{4.85}$$

Then for the tilted hyperbolic needle  $H_N(\sqrt{2}; \gamma)$  of area  $(\gamma/\sqrt{2}) \log N + O(1)$ , there exist translated copies  $\mathbf{x}_1 + H_N(\sqrt{2}; \gamma) \subset [0, M]^2$  and  $\mathbf{x}_2 + H_N(\sqrt{2}; \gamma) \subset [0, M]^2$  such that

$$|\mathcal{P} \cap (\mathbf{x}_1 + H_N(\sqrt{2}; \gamma))| > \frac{\delta\gamma}{\sqrt{2}} \log N + \delta' \log N$$

and

$$|\mathcal{P} \cap (\mathbf{x}_2 + H_N(\sqrt{2}; \gamma))| < \frac{\delta\gamma}{\sqrt{2}} \log N - \delta' \log N.$$

Note that Propositions 10 and 11 are special cases of Theorem 12, with  $\mathcal{P} = \mathbf{Z}^2$  and  $\mathcal{P} = \mathcal{A}$  respectively.

Unfortunately, the proof of Theorem 12 is rather difficult and long, and the very complicated details cover the next four sections. But the main idea is quite simple. It is basically a sophisticated application of the Riesz product.

### 4.5 Proof of Theorem 12 (I): Proving Extra Large Deviations via Riesz Product

Since the proof is long and complicated, a convenient notation here makes a big difference. It is much simpler for us to work with hyperbolic regions in the usual horizontal-vertical position instead of the tilted position. It means that, instead of working with the set  $\mathbf{Z}^2$  of lattice points in the plane and the family of tilted hyperbolic needles of a fixed quadratic irrational slope, as in the setting of Theorem 12, we rotate back. In other words, we rotate  $\mathbf{Z}^2$  by a quadratic irrational slope, and consider the family of hyperbolic needles in the usual horizontal-vertical position.

Let  $\gamma > 0$  be an arbitrary real number, and let  $N \geq 2$  be a large integer. Consider the hyperbolic region

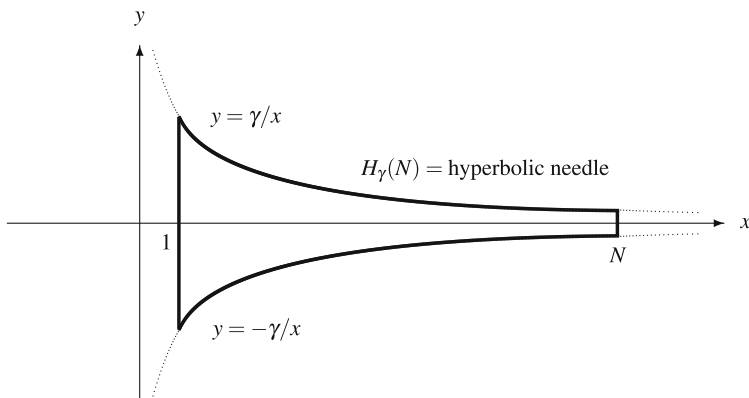


Fig. 4.2 A hyperbolic needle in usual horizontal-vertical position

$$H_\gamma(N) = \{(x, y) \in \mathbf{R}^2 : -\gamma \leq xy \leq \gamma, 1 \leq x \leq N\}; \tag{4.86}$$

see Fig. 4.2. Again we refer to  $H_\gamma(N)$  as a hyperbolic needle.

Notice that  $H_\gamma(N)$  is basically the horizontal-vertical version of the tilted hyperbolic needle  $H_N(\sqrt{2}; \gamma)$ ; see (4.56) or (4.81). To emphasize the difference between the tilted and the horizontal-vertical versions, we have made a major change in the notation, and switched the location of the parameters  $\gamma$  and  $N$ .

The area of  $H_\gamma(N)$  equals the integral

$$\text{area}(H_\gamma(N)) = 2 \int_1^N \frac{\gamma}{x} dx = 2\gamma \log N.$$

Let  $\text{rot}_\alpha \mathbf{Z}^2$  denote the rotated copy of  $\mathbf{Z}^2$  by the angle  $\theta$ , where  $\tan \theta = \alpha$  is the slope and using the origin as the fixed point of the rotation. If  $\alpha \neq 0$  is a quadratic irrational, then the continued fractions for  $\alpha$  is finally periodic. This is a well known number-theoretic fact; for example, if  $\alpha = 1/\sqrt{2}$ , then

$$\frac{1}{\sqrt{2}} = \frac{1}{1 + \frac{1}{2 + \frac{1}{2 + \frac{1}{2 + \dots}}}} = [1, 2, 2, 2, \dots] = [1, \bar{2}].$$

Periodicity implies that the continued fraction digits, formally known as the partial quotients, form a bounded sequence. It is well known that boundedness yields

$$k \|k\alpha\| \geq c_{11} = c_{11}(\alpha) > 0 \quad \text{for all integers } k \geq 1, \tag{4.87}$$

where  $c_{11} = c_{11}(\alpha) > 0$  is some positive constant depending only on  $\alpha$ , and  $\|z\|$  denotes the distance of a real number  $z$  to the nearest integer. If  $\alpha = 1/\sqrt{2}$ , then (4.87) follows from the factorization  $x^2 - 2y^2 = (x - y\sqrt{2})(x + y\sqrt{2})$ . If  $x$  and  $y$  are integers, then

$$1 \leq |x^2 - 2y^2| = |(x - y\sqrt{2})(x + y\sqrt{2})| = |x\alpha - y|\sqrt{2}|x + y\sqrt{2}|,$$

and we choose  $x = k$  and  $y$  to be the nearest integer to  $k\alpha$ . This explains why in the special case  $\alpha = 1/\sqrt{2}$  that the choice  $c_{11} = \frac{1}{4}$  in (4.87) works.

Inequality (4.87) has an important geometric interpretation, namely that there is another constant  $c_{12} = c_{12}(\alpha) > 0$ , depending only on  $\alpha$ , such that for every axes-parallel rectangle  $R$ ,

$$|\text{rot}_\alpha \mathbf{Z}^2 \cap R| \leq 1 \quad \text{whenever} \quad \text{area}(R) = c_{12}(\alpha). \tag{4.88}$$

If  $\alpha = 1/\sqrt{2}$ , then  $c_{12} = \frac{1}{5}$  is a good choice in (4.88), in view of Lemma 4.

The following statement is just a slight generalization of Theorem 12.

**Proposition 13.** *Let  $\mathcal{P}$  be a finite set of points in the square  $[0, M]^2$  with density  $\delta$ , so that the number of elements of  $\mathcal{P}$  is  $|\mathcal{P}| = \delta M^2$ . Assume further that  $\mathcal{P}$  satisfies the following rectangle property, that there is a positive constant  $c_1 = c_1(\mathcal{P}) > 0$  such that every axes-parallel rectangle of area  $c_1$  contains at most one element of the set  $\mathcal{P}$ . Let  $\delta' = \delta'(c_1, \gamma, \delta)$  be defined by (4.83) and (4.84), and assume that both  $N$  and  $M/N$  are sufficiently large and satisfy (4.85). Then for the hyperbolic needle  $H_\gamma(N)$  given by (4.86), there exist translated copies  $\mathbf{x}_1 + H_\gamma(N) \subset [0, M]^2$  and  $\mathbf{x}_2 + H_\gamma(N) \subset [0, M]^2$  such that*

$$|\mathcal{P} \cap (\mathbf{x}_1 + H_\gamma(N))| > 2\delta\gamma \log N + \delta' \log N \tag{4.89}$$

and

$$|\mathcal{P} \cap (\mathbf{x}_2 + H_\gamma(N))| < 2\delta\gamma \log N - \delta' \log N. \tag{4.90}$$

*Remarks.* (i) The term  $2\delta\gamma \log N$  in (4.89) and (4.90) represents the expectation, since the set  $\mathcal{P}$  has density  $\delta$  and the hyperbolic needle  $H_\gamma(N)$  has area  $2\gamma \log N$ . The extra terms  $\pm\delta' \log N$  show that the deviation from the expectation is proportional to the expectation, justifying the terminology *extra large deviation*.

(ii) The constant factors such as  $10^{-12}$  and  $10^{11}$  are certainly very far from best possible. Since the proof is complicated, our primary goal is to present the basic ideas in the simplest form, and we do not care too much about optimizing these constant factors.

We begin our long proof of Proposition 13.

Consider the point-counting function

$$f(\mathbf{x}) = |\mathcal{P} \cap (\mathbf{x} + H_\gamma(N))|. \tag{4.91}$$

If  $\mathbf{x} \in [0, M - N] \times [\gamma, M - \gamma]$ , then clearly

$$\mathbf{x} + H_\gamma(N) \subset [0, M]^2. \tag{4.92}$$

This explains why we choose the rectangle  $[0, M - N] \times [\gamma, M - \gamma]$  to be our underlying domain in the proof.

Let

$$\Delta(\mathbf{x}) = f(\mathbf{x}) - \delta \cdot \text{area}(H_\gamma(N)) = f(\mathbf{x}) - 2\delta\gamma \log N \quad (4.93)$$

denote the discrepancy function;  $\Delta(\mathbf{x})$  deserves its name if (4.92) holds.

In order to show that  $\Delta(\mathbf{x}) > \delta' \log N > 0$  holds for some  $\mathbf{x} = \mathbf{x}_1$ , we apply the test function method initiated by Roth [26]. The basic idea of this method is to construct a positive test function  $T(\mathbf{x}) > 0$  such that

$$\frac{1}{(M - N)(M - 2\gamma)} \int_0^{M-N} \int_\gamma^{M-\gamma} \Delta(\mathbf{x})T(\mathbf{x}) \, d\mathbf{x} > c_{13} \log N > 0, \quad (4.94)$$

and

$$\frac{1}{(M - N)(M - 2\gamma)} \int_0^{M-N} \int_\gamma^{M-\gamma} T(\mathbf{x}) \, d\mathbf{x} < c_{14}. \quad (4.95)$$

Combining (4.94) and (4.95) with the general trivial inequality

$$\int \Delta(\mathbf{x})T(\mathbf{x}) \, d\mathbf{x} \leq \max_{\mathbf{x}} \Delta(\mathbf{x}) \int T(\mathbf{x}) \, d\mathbf{x}, \quad (4.96)$$

which holds for any positive function  $T(\mathbf{x}) > 0$ , we conclude that

$$\max_{\mathbf{x}} \Delta(\mathbf{x}) > c_{15} \log N$$

with some positive constant  $c_{15} > 0$ .

Similarly, to show that  $\Delta(\mathbf{x}) < -\delta' \log N < 0$  for some  $\mathbf{x} = \mathbf{x}_2$ , we construct a positive test function  $T^*(\mathbf{x}) > 0$  such that

$$\frac{1}{(M - N)(M - 2\gamma)} \int_0^{M-N} \int_\gamma^{M-\gamma} \Delta(\mathbf{x})T^*(\mathbf{x}) \, d\mathbf{x} < -c_{16} \log N < 0, \quad (4.97)$$

and again

$$\frac{1}{(M - N)(M - 2\gamma)} \int_0^{M-N} \int_\gamma^{M-\gamma} T^*(\mathbf{x}) \, d\mathbf{x} < c_{17}. \quad (4.98)$$

Clearly (4.97) and (4.98) lead to the inequality

$$\min_{\mathbf{x}} \Delta(\mathbf{x}) < -c_{18} \log N < 0$$

with some positive constant  $c_{18} > 0$ .

Let us return to (4.94) and (4.95). We shall express the test function  $T(\mathbf{x})$  in terms of modified Rademacher functions, sometimes called Haar wavelet, and this is another idea that we borrow from Roth’s pioneering paper [26]. The benefit of working with modified Rademacher functions is that we have orthogonality and, what is more, we have *super-orthogonality*; see the key property below.

Note that Roth simply took the sum of certain modified Rademacher functions, and applied the Cauchy–Schwarz inequality instead of (4.96). For his argument, orthogonality was sufficient. It was Halász’s innovation<sup>8</sup> to express  $T(\mathbf{x})$  as a Riesz product of modified Rademacher functions; see Halász [19]. The main point is that the Riesz product takes advantage of the super-orthogonality. Here we develop an adaptation of the Roth–Halász method for hyperbolic regions.

Following the Roth–Halász approach, we shall express the test function  $T(\mathbf{x})$  as a Riesz product of modified Rademacher functions, in the form

$$T(\mathbf{x}) = \prod_{j \in \mathcal{J}} (1 + \rho R_j(\mathbf{x})), \tag{4.99}$$

where  $0 < \rho < 1$  is an appropriate constant to be specified later,  $\mathcal{J}$  is some appropriate index-set and  $R_j(\mathbf{x})$ ,  $j \in \mathcal{J}$ , are certain modified Rademacher functions to be defined below. We assume that the test function  $T(\mathbf{x})$  is zero outside the rectangle  $[0, M - N] \times [\gamma, M - \gamma]$ .

Suppose that  $10^{-2} > \eta_1 > 0$  and  $10^{-2} > \eta_2 > 0$  are small positive real numbers, to be specified later, such that

$$\frac{M - N}{\eta_1} = \frac{M - 2\gamma}{\eta_2} = 2^m, \tag{4.100}$$

where  $m \geq 1$  is an integer. Let  $j$  be an arbitrary integer in the interval  $0 \leq j \leq n$  where  $2^n \approx N$ , that is,  $n = \log_2 N + O(1)$  in binary logarithm. We decompose the rectangle  $[0, M - N] \times [\gamma, M - \gamma]$  into  $2^m \times 2^m = 4^m$  disjoint translated copies of the small rectangle

$$[0, 2^j \eta_1] \times [0, 2^{-j} \eta_2], \tag{4.101}$$

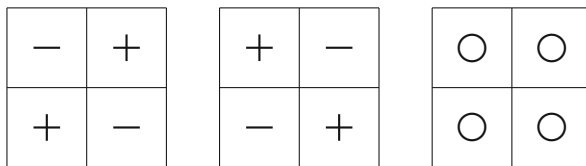
and call these congruent copies of the small rectangle (4.101)  $j$ -cells. For each of the  $4^m$   $j$ -cells, we independently choose one of the three patterns  $+-$ ,  $-+$  and  $0$ ; see Fig. 4.3.

As Fig. 4.3 shows, the pattern  $+-$  actually means a two-dimensional pattern as follows. We divide the  $j$ -cell into four congruent subrectangles, and define a step-function on the  $j$ -cell, with value  $+1$  on the upper-right and lower-left subrectangles, and value  $-1$  on the upper-left and lower-right subrectangles.

---

<sup>8</sup>Halász used this method, among many other things, to give an elegant new proof of Schmidt’s well-known discrepancy theorem; see [27].

**Fig. 4.3** The patterns  $+-$ ,  $-+$  and 0



Similarly, the pattern  $-+$  means the step-function with value  $-1$  on the upper-right and lower-left subrectangles, and value  $+1$  on the upper-left and lower-right subrectangles.

Finally, the pattern 0 means that the step-function is zero on the whole  $j$ -cell.

In the sequel, we shall simply refer to these two-dimensional patterns as  $+-$ ,  $-+$  and 0, representing the bottom rows in Fig. 4.3.

By making an independent choice of  $+-$ ,  $-+$  and 0 for each  $j$ -cell, we obtain a particular modified Rademacher function  $R_j(\mathbf{x})$  of order  $j$ , defined over the whole rectangle  $[0, M - N] \times [\gamma, M - \gamma]$ . We define  $R_j(\mathbf{x})$  to be 0 outside the rectangle  $[0, M - N] \times [\gamma, M - \gamma]$ .

Since for each of the  $4^m$   $j$ -cells there are 3 options, namely  $+-$ ,  $-+$  and 0, the total number of modified Rademacher functions  $R_j(\mathbf{x})$  of order  $j$  is  $3^{4^m}$ . Let  $\mathcal{R}(j)$  denote the family of all  $3^{4^m}$  modified Rademacher functions of order  $j$ . Note that the notation  $R_j(\mathbf{x})$  is somewhat ambiguous in the sense that it represents any element of this huge family  $\mathcal{R}(j)$ .

**Super-Orthogonality: Key Property of the Modified Rademacher Functions.** *If  $k \geq 1$  and  $0 \leq j_1 < \dots < j_k \leq n$ , then in every elementary cell of size  $2^{j_1} \eta_1 \times 2^{-j_k} \eta_2$ , the product  $R_{j_1}(\mathbf{x}) \dots R_{j_k}(\mathbf{x})$  of  $k$  modified Rademacher functions satisfies one of the three familiar patterns in Fig. 4.3.*

Note that an elementary cell of size  $2^{j_1} \eta_1 \times 2^{-j_k} \eta_2$  arises as a non-empty intersection of a  $j_1$ -cell and a  $j_k$ -cell, where  $j_1 < j_k$ . The proof of the above key property is almost trivial. It is based on the fact that for any  $k \geq 2$ , the intersection of any  $k$  cells of different orders  $j_1 < \dots < j_k$  is either empty or equal to the intersection of the  $j_1$ -cell and the  $j_k$ -cell, i.e. the intersection of the first and the last. We emphasize that in each of the 3 patterns the integral of the corresponding step-function is zero.

Since every modified Rademacher function  $R_j(\mathbf{x})$  has values  $\pm 1$  or 0, and since  $0 < \rho < 1$ , it is clear that the Riesz product (4.99) defines a positive test function  $T(\mathbf{x})$ . The index-set  $\mathcal{J}$ , a subset of  $\{0, 1, 2, \dots, n\}$ , will be specified later. Note in advance that  $\mathcal{J}$  is a large subset of  $\{0, 1, 2, \dots, n\}$ , in the sense that  $|\mathcal{J}| \geq c_{19}(n + 1)$ .

Next we check the second requirement (4.95) of the test function. Multiplying out the Riesz product (4.99), we have

$$T(\mathbf{x}) = \prod_{j \in \mathcal{J}} (1 + \rho R_j(\mathbf{x}))$$



$$\begin{aligned}
 &= 1 + \rho \sum_{j \in \mathcal{J}} R_j(\mathbf{x}) + \rho^2 \sum_{\substack{j_1 < j_2 \\ j_i \in \mathcal{J}}} R_{j_1}(\mathbf{x}) R_{j_2}(\mathbf{x}) \\
 &\quad + \rho^3 \sum_{\substack{j_1 < j_2 < j_3 \\ j_i \in \mathcal{J}}} R_{j_1}(\mathbf{x}) R_{j_2}(\mathbf{x}) R_{j_3}(\mathbf{x}) + \dots, \tag{4.102}
 \end{aligned}$$

in the form 1 plus the linear part plus the quadratic part plus the cubic part and so on. Substituting (4.102) into the left hand side of (4.95), we have

$$\begin{aligned}
 &\frac{1}{(M - N)(M - 2\gamma)} \int_0^{M-N} \int_{\gamma}^{M-\gamma} T(\mathbf{x}) \, d\mathbf{x} \\
 &= 1 + \sum_{k \geq 1} \frac{\rho^k}{(M - N)(M - 2\gamma)} \sum_{\substack{j_1 < \dots < j_k \\ j_i \in \mathcal{J}}} \int_0^{M-N} \int_{\gamma}^{M-\gamma} R_{j_1}(\mathbf{x}) \dots R_{j_k}(\mathbf{x}) \, d\mathbf{x} \\
 &= 1. \tag{4.103}
 \end{aligned}$$

The vanishing integrals in the last step occurs as a consequence of the super-orthogonality of the modified Rademacher functions. For each of 3 patterns that the integrand takes, the integral is zero. Clearly (4.103) gives (4.95) with  $c_{14} = 1$ .

Finally, we turn to requirement (4.94). The verification of this is by far the most difficult part of the proof. This is where we make the critical decision on how we choose an appropriate modified Rademacher function  $R_j(\mathbf{x})$  from amongst the huge family  $\mathcal{R}(j)$  of size  $3^{4^m}$ . We choose the best  $R_j(\mathbf{x}) \in \mathcal{R}(j)$  in order to *synchronize the trivial errors*. The synchronization argument is at the very heart of the proof. Note that if we did not synchronize the trivial errors, then they might cancel out, and we would then not be able to guarantee extra large deviation.

*The Trivial Errors and Synchronization.* By (4.91) and (4.93), the discrepancy function equals

$$\Delta(\mathbf{x}) = |\mathcal{P} \cap (\mathbf{x} + H_{\gamma}(N))| - \delta \cdot \text{area}(H_{\gamma}(N)),$$

and so we can write

$$\begin{aligned}
 &\int_0^{M-N} \int_{\gamma}^{M-\gamma} \Delta(\mathbf{x}) T(\mathbf{x}) \, d\mathbf{x} \\
 &= \int_0^{M-N} \int_{\gamma}^{M-\gamma} \left( \sum_{P_i \in \mathcal{P} \cap (\mathbf{x} + H_{\gamma}(N))} 1 - \delta \cdot \text{area}(H_{\gamma}(N)) \right) T(\mathbf{x}) \, d\mathbf{x}
 \end{aligned}$$

$$\begin{aligned}
 &= \int_0^{M-N} \int_\gamma^{M-\gamma} \left( \sum_{P_i \in \mathcal{P} \cap (\mathbf{x} + H_\gamma(N))} 1 \right) T(\mathbf{x}) \, d\mathbf{x} \\
 &\quad - (M - N)(M - 2\gamma)\delta \cdot \text{area}(H_\gamma(N)), \tag{4.104}
 \end{aligned}$$

where in the last step we have used (4.103), and where  $P_1, P_2, P_3, \dots$  denote the elements of the given point set  $\mathcal{P}$ .

Changing the order of summation and integration, we obtain

$$\int_0^{M-N} \int_\gamma^{M-\gamma} \left( \sum_{P_i \in \mathcal{P} \cap (\mathbf{x} + H_\gamma(N))} 1 \right) T(\mathbf{x}) \, d\mathbf{x} = \sum_{P_i \in \mathcal{P}} \int_{P_i - H_\gamma(N)} T(\mathbf{x}) \, d\mathbf{x}, \tag{4.105}$$

where

$$P_i - H_\gamma(N) = \{P_i - \mathbf{w} : \mathbf{w} \in H_\gamma(N)\}$$

denotes a reflected and translated copy of the hyperbolic needle  $H_\gamma(N)$ . Combining (4.104) and (4.105), we have

$$\begin{aligned}
 &\frac{1}{(M - N)(M - 2\gamma)} \int_0^{M-N} \int_\gamma^{M-\gamma} \Delta(\mathbf{x}) T(\mathbf{x}) \, d\mathbf{x} \\
 &= \sum_{P_i \in \mathcal{P}} \frac{1}{(M - N)(M - 2\gamma)} \int_{P_i - H_\gamma(N)} T(\mathbf{x}) \, d\mathbf{x} - \delta \cdot \text{area}(H_\gamma(N)). \tag{4.106}
 \end{aligned}$$

To evaluate (4.106), we return to the Riesz product (4.102). Note that the term 1 in fact denotes the characteristic function  $\chi_B$  of the rectangle  $B = [0, M - N] \times [\gamma, M - \gamma]$ , since by definition the modified Rademacher functions are all zero outside  $B$ .

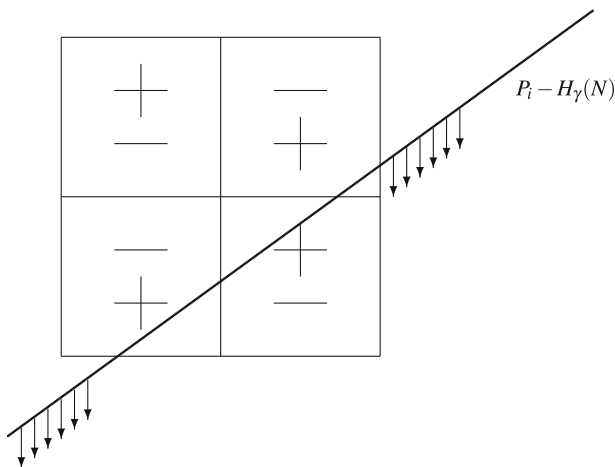
We begin with the contribution of  $1 = \chi_B$  in (4.102), and note simply that

$$\int_{P_i - H_\gamma(N)} \chi_B(\mathbf{x}) \, d\mathbf{x} = \int_{B \cap (P_i - H_\gamma(N))} d\mathbf{x} = \text{area}(B \cap (P_i - H_\gamma(N))). \tag{4.107}$$

*Geometric Ideas.* Next we study the contribution of the linear part of (4.102) in (4.106). Synchronization means that we want to make the sum

$$\sum_{P_i \in \mathcal{P}} \int_{P_i - H_\gamma(N)} R_j(\mathbf{x}) \, d\mathbf{x} \tag{4.108}$$

large and positive for every  $j \in \mathcal{J}$ , where the index-set  $\mathcal{J} \subset \{0, 1, 2, \dots, n\}$  will be specified later. We decompose the underlying rectangle  $B = [0, M - N] \times [\gamma, M - \gamma]$  into  $j$ -cells. Let  $\mathcal{C}$  be an arbitrary  $j$ -cell; it has size  $\eta_1 \eta_2$ . Consider a



**Fig. 4.4** Intersection of a  $j$ -cell with a hyperbolic arc  $P_i - H_\gamma(N)$

single term in (4.108), and restrict it to the  $j$ -cell  $\mathcal{C}$ . The geometric meaning of the integral

$$\int_{\mathcal{C} \cap (P_i - H_\gamma(N))} R_j(\mathbf{x}) \, d\mathbf{x} \tag{4.109}$$

plays a crucial role in the argument below; see Fig. 4.4.

Since the  $j$ -cell is very small, the hyperbola arc  $P_i - H_\gamma(N)$  can be approximated by its tangent line locally. This explains the tilted straight line segment in Fig. 4.4. The arrows indicate the inside of the hyperbolic needle, i.e. the arc in the picture is the upper arc of the needle.

The value of integral (4.109) depends heavily on which of the 3 patterns happens to show up in the restriction of  $R_j(\mathbf{x})$  to the  $j$ -cell  $\mathcal{C}$ . The patterns  $+ -$  and  $- +$  give two integrals whose sum is 0, whereas the pattern 0 clearly gives an integral with value 0.

How do we choose the right pattern  $+ -$ ,  $- +$  or 0 in an arbitrary  $j$ -cell  $\mathcal{C}$ ? Well, for a fixed point the choice is trivial. For every fixed point  $P_i \in \mathcal{P}$ , exactly one of the two patterns  $+ -$  and  $- +$  will make the integral (4.109) positive, unless both integrals are equal to 0. The problem is that we are dealing with a large sum

$$\sum_{P_i \in \mathcal{P}} \int_{\mathcal{C} \cap (P_i - H_\gamma(N))} R_j(\mathbf{x}) \, d\mathbf{x} \tag{4.110}$$

instead of just a single term (4.109), and we have to make (4.110) positive. The difficulty is that different points may prefer different patterns; say, for  $P_{i_1}$  the pattern

$+-$  may make the integral (4.109) positive, whereas for another point  $P_{i_2}$  the pattern  $-+$  may make the integral (4.109) positive.

To overcome this difficulty, we will apply the *Single Dominant Term Rule*, which means the following. If the sum (4.110) is dominated by a single term (4.109), then by an appropriate choice between the patterns  $+-$  and  $-+$ , we can always make this dominant term positive. We then show that the contribution from the remaining terms to (4.110) is relatively negligible. If there is no dominant term in (4.110), then we choose the pattern 0.

Of course, we have to define precisely what *domination* means. The success of the Single Dominant Term Rule is based on the fact that single term domination is quite typical: it happens very often among the  $4^m$   $j$ -cells.

What is *single term domination* in (4.110)? To explain this, we have to talk about slopes. The slope of the diagonal of a  $j$ -cell is

$$4^{-j} \eta_2 / \eta_1 \approx 4^{-j},$$

since  $\eta_1$  and  $\eta_2$  are almost equal.<sup>9</sup> Since the hyperbola is a smooth curve, the intersection of a translated and reflected hyperbolic needle  $P_i - H_\gamma(N)$  with the  $j$ -cell  $\mathcal{C}$  is almost like the intersection of  $\mathcal{C}$  with a half-plane, or the intersection of  $\mathcal{C}$  with two nearly parallel half-planes. Since half-planes have well-defined constant slopes, as an intuitive oversimplification, we shall use the terms *half-plane* and *slope* for the intersections  $\mathcal{C} \cap (P_i - H_\gamma(N))$ . Single term domination occurs if

- there is precisely one half-plane  $\mathcal{C} \cap (P_i - H_\gamma(N))$  with slope close to  $4^{-j}$  that intersects  $\mathcal{C}$ ; and
- this intersection is a *large triangle* in only one of the four subrectangles of  $\mathcal{C}$ , namely the lower right subrectangle, where the pattern is constant.

Here the intersection requirement *large triangle from the lower right subrectangle* guarantees that the integral (4.109) is far from zero, and the integral (4.109) of this dominant term is called the *trivial error*.

*An Important Consequence of the Rectangle Property.* As indicated above, single term domination means that there is exactly one half-plane  $\mathcal{C} \cap (P_i - H_\gamma(N))$  with slope close to  $4^{-j}$ . It is important to point out that we cannot have two half-planes with slopes very close to  $4^{-j}$  such that both are upper arcs. As shown Fig. 4.5, if  $\mathcal{C} \cap (P_{i_1} - H_\gamma(N))$  and  $\mathcal{C} \cap (P_{i_2} - H_\gamma(N))$  are both upper arcs with slopes very close to  $4^{-j}$ , then the two points  $P_{i_1}$  and  $P_{i_2}$  have to be in the same axes-parallel rectangle of area  $c_1$ , namely, in an axes-parallel rectangle where the slope of the

---

<sup>9</sup>We do not distinguish between positive and negative slopes. Note that the reflected hyperbolic needle  $-H_\gamma(N)$  has two long arcs: the upper arc, which is increasing, and the lower arc, which is decreasing; here the lower arc is below the upper arc. When we say that  $P_i - H_\gamma(N)$  intersects  $\mathcal{C}$ , then it always means that at least one of the two long arcs of  $P_i - H_\gamma(N)$  intersects  $\mathcal{C}$ . For example, in the *trivial error* discussed at the end of this paragraph, the intersection comes from the upper arc.

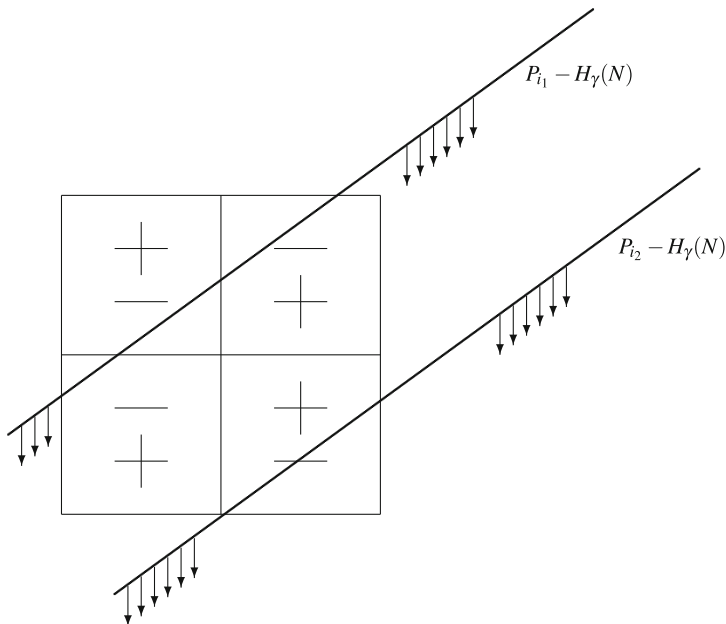
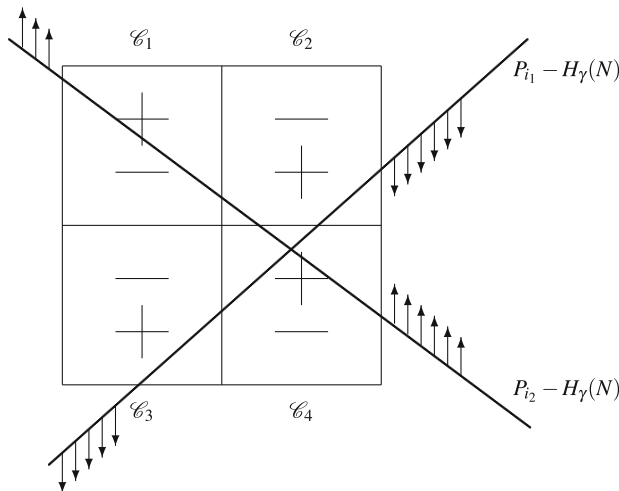


Fig. 4.5 Forbidden configuration

diagonal is close to  $4^{-j}$ . But two points in the same axes-parallel rectangle of area  $c_1$  is impossible: it contradicts the hypothesis of Proposition 13.

What can happen, however, is that we have two half-planes with slopes very close to  $4^{-j}$  such that one is an upper arc and the other one is a lower arc. For example, it can happen that  $\mathcal{C} \cap (P_{i_1} - H_\gamma(N))$  is an upper arc and  $\mathcal{C} \cap (P_{i_2} - H_\gamma(N))$  is a lower arc with both slopes<sup>10</sup> close to  $4^{-j}$ . To overcome this difficulty, we switch to a  $2 \times 2$  configuration of  $j$ -cells. More precisely, instead of working with a single  $j$ -cell  $\mathcal{C}$ , we switch to a  $2 \times 2$  configuration of four neighboring  $j$ -cells  $\mathcal{C}_1, \mathcal{C}_2, \mathcal{C}_3$  and  $\mathcal{C}_4$ , where  $\mathcal{C}_1$  is the upper left,  $\mathcal{C}_2$  is the upper right,  $\mathcal{C}_3$  is the lower left and  $\mathcal{C}_4$  is the lower right member of the  $2 \times 2$  configuration. The simple geometric idea is the following. Assume that the upper arc of  $P_{i_1} - H_\gamma(N)$  intersects both  $\mathcal{C}_2$  and  $\mathcal{C}_3$  satisfying the requirement *large triangle from the lower right subrectangle*, where the pattern is constant. Then obviously the lower arc of  $P_{i_2} - H_\gamma(N)$  cannot intersect both of  $\mathcal{C}_2$  and  $\mathcal{C}_3$ , since the slopes are close to  $4^{-j}$ . Therefore, either  $\mathcal{C}_2$  or  $\mathcal{C}_3$  will be a  $j$ -cell with single term domination. That is, we can always save at least one of the four neighboring  $j$ -cells  $\mathcal{C}_1, \mathcal{C}_2, \mathcal{C}_3$  and  $\mathcal{C}_4$ . See Fig. 4.6, where  $\mathcal{C}_3$  has single term domination.

<sup>10</sup>Again, we do not distinguish between positive and negative slopes.



**Fig. 4.6** A  $2 \times 2$  configuration of  $j$ -cells

*Choosing a Short Vertical Translation.* Next we explain how one can satisfy the intersection requirement *large triangle from the lower right subrectangle*, where the pattern is constant. This is very important, since this requirement guarantees that the dominant integral (4.109) is far from zero. First we pick an arbitrary point  $P_i \in \mathcal{P}$ . Then of course the hyperbolic needle  $P_i - H_\gamma(N)$  has a long arc such that the slope is close to  $4^{-j}$ ; *long* in fact means length of roughly  $2^j$ . Therefore, for each point  $P_i \in \mathcal{P}$ , there is a  $j$ -cell  $\mathcal{C}$  such that the intersection  $\mathcal{C} \cap (P_i - H_\gamma(N))$  has slope close to  $4^{-j}$ . Unfortunately, nothing guarantees that  $P_i - H_\gamma(N)$  intersects only one of the four subrectangles, where the pattern is constant. The solution is very simple. We apply a short vertical translation of the point set  $\mathcal{P}$ , but of course the modified Rademacher functions and the test function  $T(\mathbf{x})$  remain fixed in the rectangle  $B = [0, M - N] \times [\gamma, M - \gamma]$ . Here a *short* vertical translation means that the length of the vertical translation runs from 0 to 1. For a  $j$ -cell, a translation of length from 0 to  $2^{-j} \eta_2$  already suffices: as the point  $P_i$  moves up vertically, the intersection  $\mathcal{C} \cap (P_i - H_\gamma(N))$  changes, and has good positions where  $P_i - H_\gamma(N)$  intersects only the lower right subrectangle, where the pattern is constant, and at the same time, this intersection is a large triangle. Since the slope is close to  $4^{-j}$ , a positive constant percentage of the translations is good. If we apply translations from 0 to 1, then it will work for all  $j$ .

It follows from a standard averaging argument that there is<sup>11</sup> a vertical translation  $0 < t_0 < 1$  which is good for many pairs  $(P_i, j)$  at the same time, where  $P_i \in \mathcal{P}$  is a given point and  $j \in \{0, 1, 2, \dots, n\}$  is an order of the modified Rademacher function. Here *many* means a positive constant percentage of all pairs.

<sup>11</sup>In fact, the majority will do.

Of course, a vertical translation has a bad side effect. It causes some points to leave the underlying square  $[0, M]^2$ . However, luckily for us, it suffices to use short translations of length at most 1, so that we lose relatively few points, and only those that are close to the boundary. Note that the rectangle property in the hypothesis of Proposition 13 guarantees that there are at most  $O(M)$  points close to the boundary, which clearly is negligible compared to the number  $\delta M^2$  of points in  $\mathcal{P}$ .

*Summarizing the Vague Geometric Intuition.* A typical vertical translation of length  $0 < t_0 < 1$  has the property that for a positive constant percentage of the pairs  $(j, \mathcal{C})$ , where  $j \in \{0, 1, 2, \dots, n\}$  and  $\mathcal{C}$  is a  $j$ -cell, we have single term domination, so that<sup>12</sup>

$$\sum_{P_i \in \mathcal{D}} \int_{\mathcal{C} \cap (P_i - H_\gamma(N))} R_j(\mathbf{x}) \, d\mathbf{x} \geq \frac{1}{2} \int_{\mathcal{C} \cap (P_{i_0} - H_\gamma(N))} R_j(\mathbf{x}) \, d\mathbf{x} \geq c_{20} > 0, \quad (4.111)$$

where  $P_{i_0}$  is the dominating point, i.e. the intersection  $\mathcal{C} \cap (P_{i_0} - H_\gamma(N))$  has slope close to  $4^{-j}$ , and this intersection is a large triangle from the lower right subrectangle of  $\mathcal{C}$ , where the pattern is constant. We shall explain the missing details of (4.111) later, and give an explicit value for  $c_{20}$ .

The Single Term Domination Rule and (4.111) give

$$\begin{aligned} & \sum_{j \in \mathcal{J}} \sum_{P_i \in \mathcal{D}} \frac{1}{(M - N)(M - 2\gamma)} \int_{P_i - H_\gamma(N)} R_j(\mathbf{x}) \, d\mathbf{x} \\ & \geq c_{21} |\mathcal{J}| \geq c_{22}(n + 1) > 0. \end{aligned} \quad (4.112)$$

The geometric intuition requires that  $j \in \mathcal{J}$  satisfies an inequality like

$$\max \left\{ 1, \frac{1}{\gamma} \right\} \leq 2^j \leq \min \left\{ N, \frac{N}{\gamma} \right\}. \quad (4.113)$$

To guarantee (4.113), we choose  $\mathcal{J}$  to be the interval of integers  $j \in \{0, 1, 2, \dots, n\}$  satisfying

$$\log_2 \left( \max \left\{ 1, \frac{1}{\gamma} \right\} \right) \leq j \leq \log_2 N - \log_2 (\max\{1, \gamma\}). \quad (4.114)$$

We emphasize that this was just an intuitive proof of (4.112). We shall return to (4.111) and (4.112) later, and show how we can make the whole argument perfectly precise and explicit.

We shall complete the proof of Proposition 13 in the next three sections. Note that (4.112) is the most difficult part.

---

<sup>12</sup>Here we skip a lot of technical details!

### 4.6 Proof of Theorem 12 (II): More on the Riesz Product

*Applying Super-Orthogonality.* We next turn to the contribution of the quadratic, cubic and higher order terms of the Riesz product (4.102) to (4.106). Let  $k \geq 2$ , and let  $0 \leq j_1 < \dots < j_k \leq n$ . Suppose that  $\mathcal{C}^*$  is the non-empty intersection of  $k$  cells of orders  $j_1 < \dots < j_k$ . Then  $\mathcal{C}^*$  is an elementary cell of size  $2^{j_1} \eta_1 \times 2^{-j_k} \eta_2 = 2^{j_1-j_k} \eta_1 \eta_2$ . Super-orthogonality yields that the product  $R_{j_1}(\mathbf{x}) \dots R_{j_k}(\mathbf{x})$  of  $k$  modified Rademacher functions of the given orders, restricted to  $\mathcal{C}^*$ , equals one of the 3 patterns  $+-, -+$  or  $0$ .

Assume that the translated and reflected hyperbolic needle  $P_i - H_\gamma(N)$  intersects  $\mathcal{C}^*$ , and let slope = slope( $\mathcal{C}^* \cap (P_i - H_\gamma(N))$ ) denote the slope<sup>13</sup> of the intersection  $\mathcal{C}^* \cap (P_i - H_\gamma(N))$ . Simple geometric consideration shows that, roughly speaking, the integral

$$\frac{1}{\text{area}(\mathcal{C}^*)} \int_{\mathcal{C}^* \cap (P_i - H_\gamma(N))} R_{j_1}(\mathbf{x}) \dots R_{j_k}(\mathbf{x}) \, d\mathbf{x}$$

is negligible unless the slope of the intersection  $\mathcal{C}^* \cap (P_i - H_\gamma(N))$  is close to  $2^{-(j_1+j_k)}$ , the slope of the diagonal of  $\mathcal{C}^*$ . More precisely, we have

$$\begin{aligned} & \left| \frac{1}{\text{area}(\mathcal{C}^*)} \int_{\mathcal{C}^* \cap (P_i - H_\gamma(N))} R_{j_1}(\mathbf{x}) \dots R_{j_k}(\mathbf{x}) \, d\mathbf{x} \right| \\ & \leq \min \left\{ \frac{1}{\text{slope} \cdot 2^{j_1+j_k}}, \text{slope} \cdot 2^{j_1+j_k} \right\}. \end{aligned} \tag{4.115}$$

Note that (4.115) is a straightforward corollary of the geometry of the 3 possible patterns of  $R_{j_1}(\mathbf{x}) \dots R_{j_k}(\mathbf{x})$  in  $\mathcal{C}^*$ .

The hyperbolic needle  $H_\gamma(N)$  is bounded by the long curves  $y = \gamma/x$  and its reflection  $y = -\gamma/x$ , with  $1 \leq x \leq N$ . The slope is the derivative  $(-\gamma/x)' = \gamma x^{-2}$ . The number of elementary cells  $\mathcal{C}^*$  of size  $2^{j_1-j_k} \eta_1 \eta_2$  intersecting a fixed hyperbolic needle  $P_i - H_\gamma(N)$  is estimated from above by the simple expression

$$2 \left( \frac{2N}{2^{j_1} \eta_1} + \frac{2\gamma}{2^{-j_k} \eta_2} \right). \tag{4.116}$$

Here the factor 2 comes from the two long boundary hyperbolic curves, the first term comes from the pointed end of the hyperbolic needle, and the second term comes from the wide part of the hyperbolic needle. A more detailed explanation of (4.116) goes as follows.

---

<sup>13</sup>We do not distinguish between positive and negative slopes.



Let us start with the pointed end of the hyperbolic needle  $H_\gamma(N)$ .

*Case A.* As  $x$  runs through the interval  $N \geq x \geq \sqrt{\gamma}2^{(j_1+j_k)/2}$ , the slope of the intersection  $\mathcal{C}^* \cap (P_i - H_\gamma(N))$  is  $\gamma x^{-2}$ , which is less than  $2^{-(j_1+j_k)}$ , the slope of the diagonal of  $\mathcal{C}^*$ . It follows that in this range,  $P_i - H_\gamma(N)$  intersects fewer than

$$2 \cdot \frac{2N}{2^{j_1} \eta_1}$$

elementary cells  $\mathcal{C}^*$  of size  $2^{j_1-j_k} \eta_1 \eta_2$ , with total area not exceeding  $4\eta_2 N 2^{-j_k}$ .

*Case B.* As  $x$  runs through the interval  $\sqrt{\gamma}2^{(j_1+j_k)/2} \geq x \geq 1$ , the slope of the intersection  $\mathcal{C}^* \cap (P_i - H_\gamma(N))$  is greater than  $2^{-(j_1+j_k)}$ , the slope of the diagonal of  $\mathcal{C}^*$ . It follows that in this range,  $P_i - H_\gamma(N)$  intersects fewer than

$$2 \cdot \frac{2\gamma}{2^{-j_k} \eta_2}$$

elementary cells  $\mathcal{C}^*$  of size  $2^{j_1-j_k} \eta_1 \eta_2$ , with total area not exceeding  $4\eta_1 \gamma 2^{j_1}$ .

In Case A, we view the hyperbola  $xy = \gamma$  as  $y = \gamma/x$ . In Case B, we switch the role of the coordinate axes and view the same hyperbola as  $x = \gamma/y$ . Thus by (4.115) and (4.116), we have

$$\begin{aligned} & \left| \int_{P_i - H_\gamma(N)} R_{j_1}(\mathbf{x}) \dots R_{j_k}(\mathbf{x}) \, d\mathbf{x} \right| \\ & \leq 4\eta_2 N 2^{-j_k} \cdot \frac{2}{n} \int_{\sqrt{\gamma}2^{(j_1+j_k)/2}}^N \frac{\gamma 2^{j_1+j_k}}{x^2} \, dx + 4\eta_1 \gamma 2^{j_1} \cdot \frac{2}{\gamma} \int_{\sqrt{\gamma}2^{-(j_1+j_k)/2}}^\gamma \frac{\gamma 2^{-(j_1+j_k)}}{y^2} \, dy \\ & = 8\eta_2 2^{-j_k} \left( \sqrt{\gamma} 2^{(j_1+j_k)/2} - \frac{\gamma 2^{j_1+j_k}}{N} \right) + 8\eta_1 2^{j_1} \left( \sqrt{\gamma} 2^{-(j_1+j_k)/2} - \gamma 2^{-(j_1+j_k)} \right) \\ & \leq 8\sqrt{\gamma}(\eta_1 + \eta_2) 2^{(j_1-j_k)/2}. \end{aligned} \tag{4.117}$$

Recall that the contribution  $1 = \chi_B$  in (4.102), where  $B = [0, M - N] \times [\gamma, M - \gamma]$ . Combining (4.102), (4.106) and (4.107), we have

$$\begin{aligned} & \frac{1}{(M - N)(M - 2\gamma)} \int_0^{M-N} \int_\gamma^{M-\gamma} \Delta(\mathbf{x}) T(\mathbf{x}) \, d\mathbf{x} \\ & = \sum_{P_i \in \mathcal{P}} \frac{\text{area}(B \cap (P_i - H_\gamma(N)))}{(M - N)(M - 2\gamma)} - \delta \cdot \text{area}(H_\gamma(N)) \end{aligned}$$

$$\begin{aligned}
 & + \rho \sum_{j \in \mathcal{J}} \sum_{P_i \in \mathcal{P}} \frac{1}{(M - N)(M - 2\gamma)} \int_{P_i - H_\gamma(N)} R_j(\mathbf{x}) \, d\mathbf{x} \\
 & + \sum_{k \geq 2} \rho^k \sum_{\substack{j_1 < \dots < j_k \\ j_i \in \mathcal{J}}} \sum_{P_i \in \mathcal{P}} \frac{1}{(M - N)(M - 2\gamma)} \int_{P_i - H_\gamma(N)} R_{j_1}(\mathbf{x}) \dots R_{j_k}(\mathbf{x}) \, d\mathbf{x}.
 \end{aligned} \tag{4.118}$$

Using (4.117), it is easy to estimate the last term in (4.118). We have

$$\begin{aligned}
 & \sum_{k \geq 2} \rho^k \sum_{\substack{j_1 < \dots < j_k \\ j_i \in \mathcal{J}}} \sum_{P_i \in \mathcal{P}} \frac{1}{(M - N)(M - 2\gamma)} \left| \int_{P_i - H_\gamma(N)} R_{j_1}(\mathbf{x}) \dots R_{j_k}(\mathbf{x}) \, d\mathbf{x} \right| \\
 & \leq \sum_{k \geq 2} \rho^k \sum_{0 \leq j_1 < \dots < j_k \leq n} \sum_{P_i \in \mathcal{P}} \frac{8\sqrt{\gamma}(\eta_1 + \eta_2)2^{(j_1 - j_k)/2}}{(M - N)(M - 2\gamma)}.
 \end{aligned} \tag{4.119}$$

For convenience, let us write  $q = j_k - j_1$ . We estimate the sum

$$\sum_{k \geq 2} \rho^k \sum_{j_1=0}^{n-k+1} \sum_{q=k-1}^{n-j_1} \sum_{j_1 < j_2 < \dots < j_{k-1} < j_1 + q} 2^{-q/2}. \tag{4.120}$$

In the innermost sum in (4.120), the indices  $j_2, \dots, j_{k-1}$  can be chosen from among the  $q - 1$  numbers lying between  $j_1$  and  $j_1 + q$  in  $\binom{q-1}{k-2}$  ways. To simplify (4.120), we can let the indices  $j_1$  and  $q$  run up to  $n$ . Then we change the order of summation. Thus we have

$$\begin{aligned}
 & \sum_{k \geq 2} \rho^k \sum_{j_1=0}^{n-k+1} \sum_{q=k-1}^{n-j_1} \sum_{j_1 < j_2 < \dots < j_{k-1} < j_1 + q} 2^{-q/2} \\
 & \leq \sum_{k \geq 2} \rho^k \sum_{j_1=0}^n \sum_{q=k-1}^n \binom{q-1}{k-2} 2^{-q/2} = \sum_{j_1=0}^n \sum_{q=1}^n 2^{-q/2} \sum_{k=2}^{q+1} \rho^k \binom{q-1}{k-2}.
 \end{aligned} \tag{4.121}$$

Note that the innermost sum

$$\sum_{k=2}^{q+1} \rho^k \binom{q-1}{k-2} = \rho^2 \sum_{k=2}^{q+1} \rho^{k-2} \binom{q-1}{k-2} = \rho^2 (1 + \rho)^{q-1}.$$

It follows that if  $0 < \rho < \sqrt{2} - 1$ , then

$$\begin{aligned}
 \sum_{j_1=0}^n \sum_{q=1}^n 2^{-q/2} \sum_{k=2}^{q+1} \rho^k \binom{q-1}{k-2} &= \sum_{j_1=0}^n \sum_{q=1}^n 2^{-q/2} \rho^2 (1+\rho)^{q-1} \\
 &= \frac{(n+1)\rho^2}{\sqrt{2}} \sum_{q=1}^n \left(\frac{1+\rho}{\sqrt{2}}\right)^{q-1} \leq \frac{(n+1)\rho^2}{\sqrt{2}} \sum_{q=1}^{\infty} \left(\frac{1+\rho}{\sqrt{2}}\right)^{q-1} \\
 &= \frac{(n+1)\rho^2}{\sqrt{2}} \left(1 - \frac{1+\rho}{\sqrt{2}}\right)^{-1} = \frac{(n+1)\rho^2}{\sqrt{2}-1-\rho}.
 \end{aligned} \tag{4.122}$$

Combining (4.119)–(4.122), we obtain

**Lemma 14.** *If  $0 < \rho < \sqrt{2} - 1$ , then*

$$\begin{aligned}
 \sum_{k \geq 2} \rho^k \sum_{\substack{j_1 < \dots < j_k \\ j_i \in \mathcal{J}}} \sum_{P_i \in \mathcal{P}} \frac{1}{(M-N)(M-2\gamma)} \left| \int_{P_i - H_\gamma(N)} R_{j_1}(\mathbf{x}) \dots R_{j_k}(\mathbf{x}) \, d\mathbf{x} \right| \\
 \leq \frac{|\mathcal{P}|}{(M-N)(M-2\gamma)} \cdot 8\sqrt{\gamma}(\eta_1 + \eta_2) \cdot \frac{(n+1)\rho^2}{\sqrt{2}-1-\rho}.
 \end{aligned} \tag{4.123}$$

We return to (4.118). The contribution from the first term on the right hand side is  $o(1)$ , so that it is negligible. To see this, we recall that  $|\mathcal{P}| = \delta M^2$ , and also that  $P_i - H_\gamma(N) \subset B = [0, M-N] \times [\gamma, M-\gamma]$  for all but  $O(M)$  points  $P_i \in \mathcal{P}$ . Thus

$$\begin{aligned}
 \sum_{P_i \in \mathcal{P}} \frac{\text{area}(B \cap (P_i - H_\gamma(N)))}{(M-N)(M-2\gamma)} - \delta \cdot \text{area}(H_\gamma(N)) \\
 = \frac{\delta M^2 + O(M)}{(M-N)(M-2\gamma)} \cdot \text{area}(H_\gamma(N)) - \delta \cdot \text{area}(H_\gamma(N)) \\
 = O\left(\frac{N \log N}{M}\right) = o(1).
 \end{aligned} \tag{4.124}$$

For the second term on the right hand side of (4.118), we have the estimate (4.112). Thus combining (4.112), (4.118), (4.123) and (4.124), we obtain

$$\begin{aligned}
 \frac{1}{(M-N)(M-2\gamma)} \int_0^{M-N} \int_\gamma^{M-\gamma} \Delta(\mathbf{x}) T(\mathbf{x}) \, d\mathbf{x} \\
 \geq c_{23} \rho (n+1) - c_{24} \frac{(n+1)\rho^2}{\sqrt{2}-1-\rho} - o(1),
 \end{aligned}$$

where the constants, the first one yet unspecified, are positive and  $0 < \rho < \sqrt{2} - 1$ . By choosing a sufficiently small  $\rho$  in the range  $0 < \rho < \sqrt{2} - 1$ , we clearly have

$$\begin{aligned} & \frac{1}{(M - N)(M - 2\gamma)} \int_0^{M-N} \int_\gamma^{M-\gamma} \Delta(\mathbf{x})T(\mathbf{x}) \, d\mathbf{x} \\ & \geq c_{25}\rho(n + 1) > c_{26} \log N > 0, \end{aligned}$$

proving (4.94), and thus proving Proposition 13 in the positive direction; see (4.89). It remains to clarify the missing details in (4.111) and (4.112); see also the paragraph *Summarizing the Vague Geometric Intuition* at the end of Sect. 4.5.

*Single Term Domination: Clarifying the Technical Details.* The geometric ideas introduced in Sect. 4.5 lead to the following conclusion. At least half of the short vertical translations  $\mathcal{P} + (0, t_0)$ , where  $0 < t_0 < 1$ , of the given point set  $\mathcal{P}$  have the property that for at least 1 % of the pairs  $(j, \mathcal{C})$ , where  $j \in \{0, 1, 2, \dots, n\}$  and  $\mathcal{C}$  is a  $j$ -cell of the underlying rectangle  $B = [0, M - N] \times [\gamma, M - \gamma]$ , there is single term domination. This property includes, among other requirements to be specified later, that there is a dominating point  $P_{i_0} = P_{i_0}(j, \mathcal{C}) \in \mathcal{P}$  such that

- $\mathcal{C} \cap (P_{i_0} - H_\gamma(N))$  has slope between  $\frac{5}{6}4^{-j}$  and  $\frac{7}{6}4^{-j}$ ;
- $P_{i_0} - H_\gamma(N)$  intersects only the lower right subrectangle of  $\mathcal{C}$ , and the intersection is a large triangle, meaning that the area is at least  $\frac{1}{32}$  of the area of  $\mathcal{C}$ , that is, the area is at least  $\eta_1\eta_2/32$ .

Then, by choosing the pattern  $+-$  in the  $j$ -cell  $\mathcal{C}$ , we have

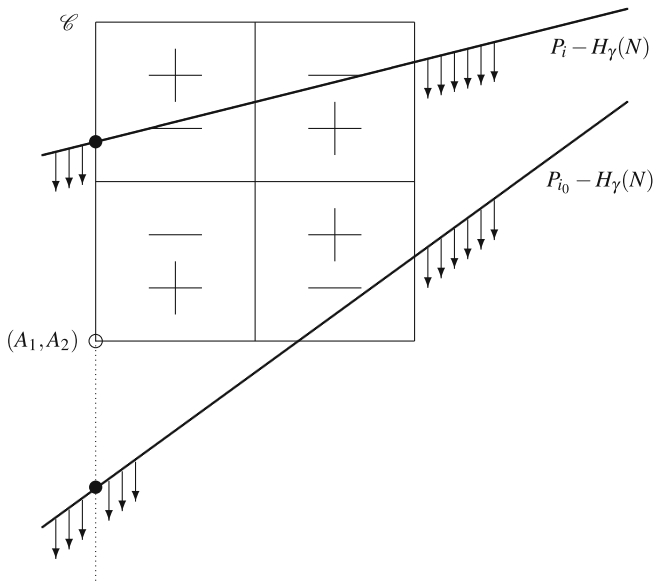
$$\int_{\mathcal{C} \cap (P_{i_0} - H_\gamma(N))} R_j(\mathbf{x}) \, d\mathbf{x} \geq \frac{\eta_1\eta_2}{32}. \tag{4.125}$$

To justify the notion *single term domination*, we shall show that for a typical pair  $(j, \mathcal{C})$ , the contribution of the remaining points  $P_i \in \mathcal{P}$ , with  $i \neq i_0$ , in the  $j$ -cell  $\mathcal{C}$  is negligible, in the sense that

$$\left| \sum_{\substack{P_i \in \mathcal{P} \\ i \neq i_0}} \int_{\mathcal{C} \cap (P_i - H_\gamma(N))} R_j(\mathbf{x}) \, d\mathbf{x} \right| \leq \frac{\eta_1\eta_2}{40}. \tag{4.126}$$

To prove (4.126), let  $P_i \neq P_{i_0}$  be another point in  $\mathcal{P}$  such that  $P_i - H_\gamma(N)$  intersects  $\mathcal{C}$ , i.e. the upper or lower arc of the boundary of the hyperbolic needle  $P_i - H_\gamma(N)$  intersects the  $j$ -cell  $\mathcal{C}$ . We are going to distinguish four cases, depending on the type of the intersection of  $P_i - H_\gamma(N)$  with  $\mathcal{C}$ , corresponding to upper or lower arc, and close to horizontal or close to vertical, relative to the diagonals of  $\mathcal{C}$ .

*Case 15.* The upper arc of  $P_i - H_\gamma(N)$  intersects  $\mathcal{C}$ , and the slope is less than the slope of the dominant needle  $P_{i_0} - H_\gamma(N)$ ; see Fig. 4.7.



**Fig. 4.7** Upper arc of  $P_i - H_\gamma(N)$  intersects  $\mathcal{C}$ , with slope less than slope of  $P_{i_0} - H_\gamma(N)$

Let  $P_{i_0} = (a_{i_0}, b_{i_0})$  and  $P_i = (a_i, b_i)$  denote the coordinates of the two points in question. By the hypothesis of Case 1, we have  $a_i > a_{i_0}$ . Write

$$h = h_i = a_i - a_{i_0} > 0 \quad \text{and} \quad v = v_i = b_i - b_{i_0},$$

where of course  $h$  denotes horizontal and  $v$  denotes vertical. The rectangle property guarantees that  $h|v| \geq c_1 > 0$ .

Let  $(A_1, A_2)$  denote the coordinates of the lower left vertex of the  $j$ -cell  $\mathcal{C}$ . The intersection of the line  $x = A_1$  with the upper arcs of  $P_{i_0} - H_\gamma(N)$  and  $P_i - H_\gamma(N)$  give two points, and the hypothesis of Case 1 implies that these intersection points are close to each other. More precisely, with  $x = 1 + a_i - A_1$ , where  $a_{i_0} - A_1 > 0$  and the additional term 1 comes from the fact that the hyperbolic needle  $H_\gamma(N)$  begins at  $x = 1$ , we have the upper bound

$$\left| \left( b_{i_0} + \frac{\gamma}{x} \right) - \left( b_i + \frac{\gamma}{x+h} \right) \right| < 2 \cdot 2^{-j} \eta_2. \tag{4.127}$$

Since  $b_i - b_{i_0} = v$ , we can rewrite (4.127) in the form

$$\left| \left( \frac{\gamma}{x} - \frac{\gamma}{x+h} \right) - v \right| = \left| \frac{\gamma h}{x(x+h)} - v \right| < 2^{-j+1} \eta_2. \tag{4.128}$$

On the other hand, we know that the slope of the upper arc of  $\mathcal{C} \cap (P_{i_0} - H_\gamma(N))$  satisfies the inequality

$$\frac{5}{6}4^{-j} \leq \frac{\gamma}{x^2} \leq \frac{7}{6}4^{-j}. \tag{4.129}$$

We claim that if  $\eta_1$ , and so also  $\eta_2$ , is a small constant, then the upper arc of  $P_{i_0} - H_\gamma(N)$  intersects a large number of  $j$ -cells different from  $\mathcal{C}$  such that the slope is still almost equal to  $4^{-j}$ . Indeed, the horizontal size of  $\mathcal{C}$  is  $2^j \eta_1$  and, assuming that (4.129) holds, the inequality

$$\frac{5}{6}4^{-j} \leq \frac{\gamma}{(x + \ell 2^j \eta_1)^2} \leq \frac{7}{6}4^{-j} \tag{4.130}$$

has constant times  $1/\eta_1$  consecutive integer solutions in  $\ell$ . If  $\eta_1 > 0$  is small, then of course  $1/\eta_1$  is large, justifying our claim.

Returning to (4.128) and (4.129), and then substituting  $x$  by  $x + \ell 2^j \eta_1$ , we have the respective inequalities

$$\left| \frac{\gamma h}{(x + \ell 2^j \eta_1)(x + \ell 2^j \eta_1 + h)} - v \right| < 2^{-j+1} \eta_2 \tag{4.131}$$

and (4.130). If (4.129) holds, then there are at least  $\sqrt{\gamma}/10\eta_1$  consecutive integer solutions  $\ell$  of (4.130).

The basic idea is the following. If  $\ell$  runs through these integer solutions of (4.130) while  $\gamma, x, h$  and  $v$  remain fixed, then the function

$$\frac{\gamma h}{(x + \ell 2^j \eta_1)(x + \ell 2^j \eta_1 + h)}, \tag{4.132}$$

as a function of  $\ell$ , has substantially different values, and we expect only very few of them to be very close to a fixed  $v$  in the quantitative sense of (4.131). Of course, here we assume that  $\eta_2$  is small.

Next we work out the details of this intuition. We begin by noting that (4.130) implies

$$\sqrt{\frac{6\gamma}{5}} 2^j \geq x + \ell 2^j \eta_1 \geq \sqrt{\frac{6\gamma}{7}} 2^j. \tag{4.133}$$

Using this in (4.132), we have the good approximation

$$\frac{\gamma h}{(x + \ell 2^j \eta_1)(x + \ell 2^j \eta_1 + h)} \approx \frac{\gamma h}{\sqrt{\gamma} 2^j (\sqrt{\gamma} 2^j + h)} = \frac{h}{2^j (2^j + h/\sqrt{\gamma})}. \tag{4.134}$$

We now distinguish two cases. First assume that  $0 < h \leq \sqrt{c_1}2^{j-1}$ , where  $c_1 > 0$  is the positive constant in the rectangle property. Then the rectangle property yields

$$|v| \geq \frac{c_1}{h} \geq \frac{c_1}{\sqrt{c_1}2^{j-1}} = 2\sqrt{c_1}2^{-j} \quad (4.135)$$

and

$$\frac{h}{2^j(2^j + h/\sqrt{\gamma})} < \frac{h}{2^j 2^j} \leq \frac{\sqrt{c_1}}{2} 2^{-j}. \quad (4.136)$$

The assumption

$$\eta_2 < \frac{\sqrt{c_1}}{2}, \quad (4.137)$$

together with (4.134)–(4.136), implies that (4.131) has no solution.

We can assume, therefore, that the lower bound

$$h > \sqrt{c_1}2^{j-1} \quad (4.138)$$

holds. Now we go back to the basic idea. We claim that if we switch  $\ell$  to  $\ell + 1$  in the function (4.132), then its value changes by at least as much as

$$\frac{\eta_1 2^{-j-2}}{1 + \sqrt{\gamma}/c_1}. \quad (4.139)$$

Indeed, by (4.133), we have

$$\frac{\gamma h}{(x + \ell 2^j \eta_1)(x + \ell 2^j \eta_1 + h)} \approx \frac{1}{\sqrt{\gamma} 2^j + 2^j \eta_1} \cdot \frac{\gamma h}{\sqrt{\gamma} 2^j + h}. \quad (4.140)$$

We also have the routine estimate

$$\begin{aligned} \frac{1}{\sqrt{\gamma} 2^j} - \frac{1}{\sqrt{\gamma} 2^j + 2^j \eta_1} &= \frac{1}{\sqrt{\gamma} 2^j} \left( 1 - \frac{1}{1 + \eta_1/\sqrt{\gamma}} \right) \\ &= \frac{1}{\sqrt{\gamma} 2^j} \left( \frac{\eta_1}{\sqrt{\gamma}} - \left( \frac{\eta_1}{\sqrt{\gamma}} \right)^2 + \left( \frac{\eta_1}{\sqrt{\gamma}} \right)^3 \mp \dots \right) \approx \frac{\eta_1}{\gamma 2^j}. \end{aligned} \quad (4.141)$$

Furthermore, by (4.138), we have

$$\frac{\gamma h}{\sqrt{\gamma} 2^j + h} > \frac{\gamma}{2\sqrt{\gamma}/c_1 + 1}. \quad (4.142)$$

Then the error estimate (4.139) follows on combining (4.140)–(4.142).

Let us return to (4.132) and (4.139), and apply them in (4.131). We deduce that among the constant times  $1/\eta_1$  consecutive integer values of  $\ell$  satisfying (4.130), there are only constant times  $(1 + \sqrt{\gamma/c_1})$  that will satisfy (4.131). More explicitly, it is safe to say that

$$\text{at most } 10 \left( 1 + \sqrt{\frac{\gamma}{c_1}} \right) \text{ values of } \ell \text{ will satisfy both (4.130) and (4.131).} \tag{4.143}$$

The next step is

*A Combination of the Rectangle Property and the Pigeonhole Principle.* We recall (4.138), that  $h > \sqrt{c_1}2^j$ . Consider the power-of-two type decomposition

$$2^{r-1}\sqrt{c_1}2^j < h \leq 2^r\sqrt{c_1}2^j, \quad r = 0, 1, 2, \dots \tag{4.144}$$

We claim that for a fixed point  $P_{i_0} = (a_{i_0}, b_{i_0}) \in \mathcal{P}$  and for a fixed integer  $r \geq 0$ , there are at most

$$10\sqrt{\frac{\gamma}{c_1}}2^r \tag{4.145}$$

other points  $P_i = (a_i, b_i) \in \mathcal{P}$ , with  $P_i \neq P_{i_0}$ , such that  $h = h_i = a_i - a_{i_0} > 0$  and  $v = v_i = b_i - b_{i_0}$  satisfy (4.131), thus implicitly (4.130) also, and (4.144).

To establish the bound (4.145), first note that if  $h = h_i$  satisfies (4.144), then by (4.134) and (4.144), we have

$$\begin{aligned} \frac{\gamma h}{(x + \ell 2^j \eta_1)(x + \ell 2^j \eta_1 + h)} &\approx \frac{h}{2^j(2^j + h/\sqrt{\gamma})} \\ &\approx \frac{2^r\sqrt{c_1}2^j}{2^j(2^j + 2^r\sqrt{c_1}2^j/\sqrt{\gamma})} = \frac{2^{-j}}{1/\sqrt{\gamma} + 2^{-r}/\sqrt{c_1}}, \end{aligned}$$

so that a solution of (4.131) gives the approximation

$$v = v_i \approx 2^{-j} \left( \frac{1}{1/\sqrt{\gamma} + 2^{-r}/\sqrt{c_1}} \pm 2\eta_2 \right). \tag{4.146}$$

Assuming

$$\eta_2 < \frac{1}{8(1/\sqrt{\gamma} + 1/\sqrt{c_1})}, \tag{4.147}$$

then (4.146) yields the good approximation

$$v = v_i \approx \frac{2^{-j}}{1/\sqrt{\gamma} + 2^{-r}/\sqrt{c_1}}. \tag{4.148}$$



Suppose on the contrary that there are more than (4.145) other points  $P_i = (a_i, b_i) \in \mathcal{P}$ , with  $P_i \neq P_{i_0}$ , such that  $h = h_i = a_i - a_{i_0} > 0$  and  $v = v_i = b_i - b_{i_0}$  satisfy (4.131), thus implicitly (4.130) also, and (4.144). Then by the Pigeonhole Principle and (4.148), there must exist two points  $P_{i_1}, P_{i_2} \in \mathcal{P}$ , with  $i_1 \neq i_2$ , such that

$$v_{i_1} \approx \frac{2^{-j}}{1/\sqrt{\gamma} + 2^{-r}/\sqrt{c_1}} \approx v_{i_2} \quad \text{and} \quad |h_{i_1} - h_{i_2}| \leq \frac{2^r \sqrt{c_1} 2^j}{10\sqrt{\gamma}/c_1 2^r} = \frac{c_1 2^j}{10\sqrt{\gamma}}.$$

Since the product

$$\frac{2^{-j}}{1/\sqrt{\gamma} + 2^{-r}/\sqrt{c_1}} \cdot \frac{c_1 2^j}{\sqrt{\gamma}} = \frac{c_1}{1 + 2^{-r} \sqrt{\gamma}/c_1} < c_1,$$

we conclude that there exists an axes-parallel rectangle of area less than  $c_1$  and which contains at least two points of  $\mathcal{P}$ , namely  $P_{i_1}$  and  $P_{i_2}$ . This contradicts the rectangle property, and establishes the bound (4.145).

If  $h = h_i$  falls into the interval (4.144), then

$$\text{slope}(\mathcal{C} \cap (P_i - H_\gamma(N))) = \frac{\gamma}{(x+h)^2} \leq \frac{\gamma}{h^2} \approx \frac{\gamma}{c_1 4^r} \cdot 4^{-j}, \tag{4.149}$$

where  $4^{-j}$  almost equals the slope of the diagonals of the  $j$ -cell  $\mathcal{C}$ . By (4.149), we have

$$\frac{1}{\text{area}(\mathcal{C})} \left| \int_{\mathcal{C} \cap (P_i - H_\gamma(N))} R_j(\mathbf{x}) \, d\mathbf{x} \right| \leq \frac{10\gamma}{c_1 4^r}. \tag{4.150}$$

Furthermore, (4.150) holds for all  $j$ -cells  $\mathcal{C}$  satisfying

$$\frac{5}{6} 4^{-j} \leq \text{slope}(\mathcal{C} \cap (P_{i_0} - H_\gamma(N))) \leq \frac{7}{6} 4^{-j}. \tag{4.151}$$

Let us return now to (4.126). Combining (4.143)–(4.145) and (4.150), we have

$$\begin{aligned} & \sum_{\substack{P_i \in \mathcal{P} \\ i \neq i_0}} \sum_{\mathcal{C}} \frac{1}{\text{area}(\mathcal{C})} \left| \int_{\mathcal{C} \cap (P_i - H_\gamma(N))} R_j(\mathbf{x}) \, d\mathbf{x} \right| \leq \sum_{r \geq 0} 10 \left(1 + \sqrt{\frac{\gamma}{c_1}}\right) 10 \sqrt{\frac{\gamma}{c_1}} 2^r \frac{10\gamma}{c_1 4^r} \\ & \text{Case 1} \\ & = 1000 \left( \left(\frac{\gamma}{c_1}\right)^{3/2} + \left(\frac{\gamma}{c_1}\right)^2 \right) \sum_{r \geq 0} 2^{-r} = 2000 \left( \left(\frac{\gamma}{c_1}\right)^{3/2} + \left(\frac{\gamma}{c_1}\right)^2 \right). \end{aligned} \tag{4.152}$$

Since there are at least  $\gamma/10\eta_1$  consecutive integer solutions  $\ell$  of (4.130), assuming that (4.129) holds, we have

$$\sum_{\mathcal{C}} 1 \geq \frac{\sqrt{\gamma}}{10\eta_1}. \tag{4.153}$$

(4.151)

Recall that in order to prove (4.126), we distinguish four cases. Inequalities (4.152) and (4.153) complete Case 1. The remaining three cases will be discussed in the next section. Note that these cases are quite similar to Case 1, but there are some annoying differences in the minor details. We shall complete the proof of Proposition 13 in Sect. 4.8.

### 4.7 Proof of Theorem 12 (III): Completing the Case Study

Let us return to (4.125) and (4.126). Again we assume that there is a dominating point  $P_{i_0} = P_{i_0}(j, \mathcal{C}) \in \mathcal{P}$  such that

- $\mathcal{C} \cap (P_{i_0} - H_\gamma(N))$  has slope between  $\frac{5}{6}4^{-j}$  and  $\frac{7}{6}4^{-j}$ ;
- $P_{i_0} - H_\gamma(N)$  intersects only the lower right subrectangle of  $\mathcal{C}$ , and the intersection is a large triangle, meaning that the area is at least  $\frac{1}{32}$  of the area of  $\mathcal{C}$ , that is, the area is at least  $\eta_1\eta_2/32$ .

Again let  $P_i \neq P_{i_0}$  be another point in  $\mathcal{P}$  such that  $P_i - H_\gamma(N)$  intersects  $\mathcal{C}$ , i.e. the upper or lower arc of the boundary of the hyperbolic needle  $P_i - H_\gamma(N)$  intersects the  $j$ -cell  $\mathcal{C}$ . We now discuss the second case, which is quite similar to the first case. Roughly speaking, we switch the roles of the horizontal and the vertical.

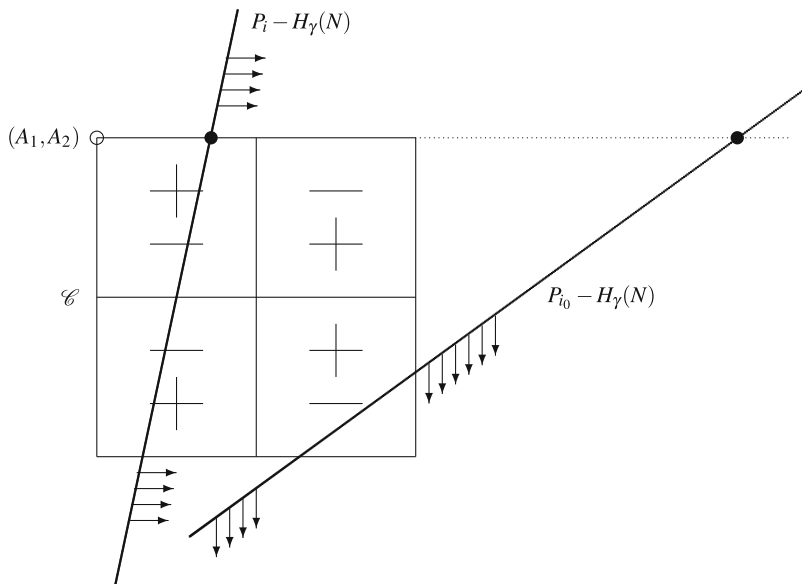
*Case 16.* The upper arc of  $P_i - H_\gamma(N)$  intersects  $\mathcal{C}$ , and the slope is greater than the slope of the dominant needle  $P_{i_0} - H_\gamma(N)$ ; see Fig. 4.8.

Let  $P_{i_0} = (a_{i_0}, b_{i_0})$  and  $P_i = (a_i, b_i)$  denote the coordinates of the two points in question. By the hypothesis of Case 2, we have  $a_{i_0} > a_i$ . Write

$$h = h_i = a_{i_0} - a_i > 0 \quad \text{and} \quad v = v_i = b_{i_0} - b_i,$$

where again  $h$  denotes horizontal and  $v$  denotes vertical. The rectangle property guarantees that  $h|v| \geq c_1 > 0$ .

Let  $(A_1, A_2)$  denote the coordinates of the upper left vertex of the  $j$ -cell  $\mathcal{C}$ . The intersection of the line  $y = A_2$  with the upper arcs of  $P_{i_0} - H_\gamma(N)$  and  $P_i - H_\gamma(N)$  give two points, and the hypothesis of Case 16 implies that these intersection points are close to each other. More precisely, with  $y = A_2 - b_{i_0}$ , we have the upper bound



**Fig. 4.8** Upper arc of  $P_i - H_\gamma(N)$  intersects  $\mathcal{C}$ , with slope greater than slope of  $P_{i_0} - H_\gamma(N)$

$$\left| \left( a_i - \frac{\gamma}{y+v} \right) - \left( a_{i_0} - \frac{\gamma}{y} \right) \right| < 2 \cdot 2^j \eta_1. \tag{4.154}$$

Since  $a_{i_0} - a_i = h > 0$ , we can rewrite (4.154) in the form

$$\left| \left( \frac{\gamma}{y} - \frac{\gamma}{y+v} \right) - h \right| = \left| \frac{\gamma v}{y(y+v)} - h \right| < 2^{j+1} \eta_1. \tag{4.155}$$

We emphasize that  $y + v > 0$ , otherwise

$$0 \geq y + v = (A_2 - b_{i_0}) + (b_{i_0} - b_i) = A_2 - b_i,$$

so that  $b_i \geq A_2$ , which means that the whole upper arc of  $P_i - H_\gamma(N)$  is above the  $j$ -cell  $\mathcal{C}$ . But this is impossible, since in Case 2 we assume that the upper arc of  $P_i - H_\gamma(N)$  intersects  $\mathcal{C}$ .

Since we switch the roles of the horizontal and the vertical, we focus on the reciprocal of the slope. We know that the reciprocal of the slope of the upper arc of  $\mathcal{C} \cap (P_{i_0} - H_\gamma(N))$  satisfies the inequality

$$\frac{6}{7} 4^j \leq \frac{\gamma}{y^2} \leq \frac{6}{5} 4^j. \tag{4.156}$$

We claim that if  $\eta_2$ , and so also  $\eta_1$ , is a small constant, then the upper arc of  $P_{i_0} - H_\gamma(N)$  intersects a large number of  $j$ -cells different from  $\mathcal{C}$  such that the reciprocal of the slope is still almost equal to  $4^j$ . Indeed, the vertical size of  $\mathcal{C}$  is  $2^{-j}\eta_2$  and, assuming that (4.156) holds, the inequality

$$\frac{6}{7}4^j \leq \frac{\gamma}{(y + \ell 2^{-j}\eta_2)^2} \leq \frac{6}{5}4^j \tag{4.157}$$

has constant times  $1/\eta_2$  consecutive integer solutions in  $\ell$ . If  $\eta_2 > 0$  is small, then of course  $1/\eta_2$  is large, justifying our claim.

Returning to (4.155) and (4.156), and then substituting  $y$  by  $y + \ell 2^{-j}\eta_2$ , we have the respective inequalities

$$\left| \frac{\gamma v}{(y + \ell 2^{-j}\eta_2)(y + \ell 2^{-j}\eta_2 + v)} - h \right| < 2^{j+1}\eta_1 \tag{4.158}$$

and (4.157). If (4.156) holds, then there are at least  $\sqrt{\gamma}/10\eta_2$  consecutive integer solutions  $\ell$  of (4.157).

The basic idea is the same as in Case 15. If  $\ell$  runs through these integer solutions of (4.157) while  $\gamma, y, h$  and  $v$  remain fixed, then the function

$$\frac{\gamma v}{(y + \ell 2^{-j}\eta_2)(y + \ell 2^{-j}\eta_2 + v)}, \tag{4.159}$$

as a function of  $\ell$ , has substantially different values, and we expect only very few of them to be very close to a fixed  $h$  in the quantitative sense of (4.158). Of course, here we assume that  $\eta_1$  is small.

Next we work out the details of this intuition. We begin by noting that (4.157) implies

$$\sqrt{\frac{6\gamma}{7}}2^{-j} \leq y + \ell 2^{-j}\eta_2 \leq \sqrt{\frac{6\gamma}{5}}2^{-j}. \tag{4.160}$$

Using this in (4.159), we have the good approximation

$$\frac{\gamma v}{(y + \ell 2^{-j}\eta_2)(y + \ell 2^{-j}\eta_2 + v)} \approx \frac{\gamma v}{\sqrt{\gamma}2^{-j}(\sqrt{\gamma}2^{-j} + v)} = \frac{v}{2^{-j}(2^{-j} + v\sqrt{\gamma})}. \tag{4.161}$$

We now distinguish three cases. First assume that  $v < 0$ . Since  $y + v > 0$ , we have  $y^{-1} < (y + v)^{-1}$ , and so by (4.158), we have

$$2^{j+1}\eta_1 > |h| = h. \tag{4.162}$$

Combining (4.162) with the rectangle property, we deduce that

$$|v| \geq \frac{c_1}{h} > \frac{c_1}{2\eta_1} 2^{-j}. \quad (4.163)$$

Substituting (4.163) into (4.161), and assuming that

$$\eta_1 < \frac{c_1}{2\sqrt{\gamma}}, \quad (4.164)$$

we have

$$\frac{v}{2^{-j}(2^{-j} + v/\sqrt{\gamma})} = \frac{|v|}{2^{-j}(v/\sqrt{\gamma} - 2^{-j})} = \frac{2^j}{1/\sqrt{\gamma} - 2^{-j}/|v|} > \sqrt{\gamma} 2^j. \quad (4.165)$$

Combining (4.158), (4.161)–(4.163) and (4.165), we conclude that

$$2^{j+1}\eta_1 > h > \frac{1}{2}\sqrt{\gamma} 2^j - 2^{j+1}\eta_1,$$

which is an obvious contradiction if

$$\eta_1 < \frac{\sqrt{\gamma}}{8}. \quad (4.166)$$

This proves that  $v > 0$ .

Next assume that  $0 < v \leq \sqrt{c_1} 2^{-j-1}$ , where  $c_1 > 0$  is the positive constant in the rectangle property. Then the rectangle property yields

$$h \geq \frac{c_1}{v} \geq \frac{c_1}{\sqrt{c_1} 2^{-j-1}} = 2\sqrt{c_1} 2^j \quad (4.167)$$

and

$$\frac{v}{2^{-j}(2^{-j} + v/\sqrt{\gamma})} < \frac{v}{2^{-j} 2^{-j}} \leq \frac{\sqrt{c_1}}{2} 2^j. \quad (4.168)$$

The assumption

$$\eta_1 < \frac{\sqrt{c_1}}{2}, \quad (4.169)$$

together with (4.161), (4.167) and (4.168), implies that (4.158) has no solution.

We can assume, therefore, that the lower bound

$$v > \sqrt{c_1} 2^{-j-1} \quad (4.170)$$

holds. Now we go back to the basic idea. We claim that if we switch  $\ell$  to  $\ell + 1$  in the function (4.159), then its value changes by at least as much as

$$\frac{\eta_1 2^{j-2}}{1 + \sqrt{\gamma/c_1}}. \tag{4.171}$$

Indeed, by (4.160), we have

$$\frac{\gamma v}{(y + \ell 2^{-j} \eta_2)(y + \ell 2^{-j} \eta_2 + v)} \approx \frac{1}{\sqrt{\gamma} 2^{-j}} \cdot \frac{\gamma v}{\sqrt{\gamma} 2^{-j} + v}. \tag{4.172}$$

We also have the routine estimate

$$\begin{aligned} \frac{1}{\sqrt{\gamma} 2^{-j}} - \frac{1}{\sqrt{\gamma} 2^{-j} + 2^{-j} \eta_2} &= \frac{1}{\sqrt{\gamma} 2^{-j}} \left( 1 - \frac{1}{1 + \eta_2/\sqrt{\gamma}} \right) \\ &= \frac{1}{\sqrt{\gamma} 2^{-j}} \left( \frac{\eta_2}{\sqrt{\gamma}} - \left( \frac{\eta_2}{\sqrt{\gamma}} \right)^2 + \left( \frac{\eta_2}{\sqrt{\gamma}} \right)^3 \mp \dots \right) \approx \frac{\eta_2 2^j}{\gamma}. \end{aligned} \tag{4.173}$$

Furthermore, by (4.170), we have

$$\frac{\gamma v}{\sqrt{\gamma} 2^{-j} + v} > \frac{\gamma}{2\sqrt{\gamma/c_1} + 1}. \tag{4.174}$$

The error estimate (4.171) follows on combining (4.172)–(4.174).

Let us return to (4.159) and (4.171), and apply them in (4.158). We deduce that among the constant times  $1/\eta_2$  consecutive integer values of  $\ell$  satisfying (4.157), there are only constant times  $(1 + \sqrt{\gamma/c_1})$  that will satisfy (4.158). More explicitly, it is safe to say that

$$\text{at most } 10 \left( 1 + \sqrt{\frac{\gamma}{c_1}} \right) \text{ values of } \ell \text{ will satisfy both (4.157) and (4.158).} \tag{4.175}$$

As in Case 15, the next step is

*A Combination of the Rectangle Property and the Pigeonhole Principle.* We recall (4.170), that  $v > \sqrt{c_1} 2^{-j-1}$ . Consider the power-of-two type decomposition

$$2^{r-1} \sqrt{c_1} 2^{-j} < v \leq 2^r \sqrt{c_1} 2^{-j}, \quad r = 0, 1, 2, \dots \tag{4.176}$$

We claim that for a fixed point  $P_{i_0} = (a_{i_0}, b_{i_0}) \in \mathcal{P}$  and for a fixed integer  $r \geq 0$ , there are at most

$$10 \sqrt{\frac{\gamma}{c_1}} 2^r \tag{4.177}$$

other points  $P_i = (a_i, b_i) \in \mathcal{P}$ , with  $P_i \neq P_{i_0}$ , such that  $h = h_i = a_{i_0} - a_i > 0$  and  $v = v_i = b_{i_0} - b_i > 0$  satisfy (4.158), thus implicitly (4.157) also, and (4.176).

To establish the bound (4.177), first note that if  $v = v_i$  satisfies (4.176), then by (4.161) and (4.176), we have

$$\begin{aligned} \frac{\gamma v}{(y + \ell 2^{-j} \eta_2)(y + \ell 2^{-j} \eta_2 + v)} &\approx \frac{v}{2^{-j}(2^{-j} + v/\sqrt{\gamma})} \\ &\approx \frac{2^r \sqrt{c_1} 2^{-j}}{2^{-j}(2^{-j} + 2^r \sqrt{c_1} 2^{-j} / \sqrt{\gamma})} = \frac{2^j}{1/\sqrt{\gamma} + 2^{-r}/\sqrt{c_1}}, \end{aligned}$$

so that a solution of (4.158) gives the approximation

$$h = h_i \approx 2^j \left( \frac{1}{1/\sqrt{\gamma} + 2^{-r}/\sqrt{c_1}} \pm 2\eta_1 \right). \tag{4.178}$$

Assuming

$$\eta_1 < \frac{1}{8(1/\sqrt{\gamma} + 1/\sqrt{c_1})}, \tag{4.179}$$

then (4.178) yields the good approximation

$$h = h_i \approx \frac{2^j}{1/\sqrt{\gamma} + 2^{-r}/\sqrt{c_1}}. \tag{4.180}$$

Suppose on the contrary that there are more than (4.177) other points  $P_i = (a_i, b_i) \in \mathcal{P}$ , with  $P_i \neq P_{i_0}$ , such that  $h = h_i = a_{i_0} - a_i > 0$  and  $v = v_i = b_{i_0} - b_i > 0$  satisfy (4.158), thus implicitly (4.157) also, and (4.176). Then by the Pigeonhole Principle and (4.180), there must exist two points  $P_{i_1}, P_{i_2} \in \mathcal{P}$ , with  $i_1 \neq i_2$ , such that

$$h_{i_1} \approx \frac{2^j}{1/\sqrt{\gamma} + 2^{-r}/\sqrt{c_1}} \approx h_{i_2} \quad \text{and} \quad |v_{i_1} - v_{i_2}| \leq \frac{2^r \sqrt{c_1} 2^{-j}}{10\sqrt{\gamma}/c_1 2^r} = \frac{c_1 2^{-j}}{10\sqrt{\gamma}}.$$

Since the product

$$\frac{2^j}{1/\sqrt{\gamma} + 2^{-r}/\sqrt{c_1}} \cdot \frac{c_1 2^{-j}}{\sqrt{\gamma}} = \frac{c_1}{1 + 2^{-r} \sqrt{\gamma}/c_1} < c_1,$$

we conclude that there exists an axes-parallel rectangle of area less than  $c_1$  and which contains at least two points of  $\mathcal{P}$ , namely  $P_{i_1}$  and  $P_{i_2}$ . This contradicts the rectangle property, and establishes the bound (4.177).

If  $v = v_i$  falls into the interval (4.176), then

$$\frac{1}{\text{slope}(\mathcal{C} \cap (P_i - H_\gamma(N)))} = \frac{\gamma}{(y + v)^2} \leq \frac{\gamma}{v^2} \approx \frac{\gamma}{c_1 4^r} \cdot 4^j, \tag{4.181}$$

where  $4^j$  almost equals the reciprocal of the slope of the diagonals of the  $j$ -cell  $\mathcal{C}$ . By (4.181), we have

$$\frac{1}{\text{area}(\mathcal{C})} \left| \int_{\mathcal{C} \cap (P_i - H_\gamma(N))} R_j(\mathbf{x}) \, d\mathbf{x} \right| \leq \frac{10\gamma}{c_1 4^r}. \tag{4.182}$$

Furthermore, (4.182) holds for all  $j$ -cells  $\mathcal{C}$  satisfying (4.151). Let us return now to (4.126). Combining (4.175)–(4.177) and (4.182), we have

$$\begin{aligned} \sum_{\substack{P_i \in \mathcal{D} \\ i \neq i_0}} \sum_{\substack{\mathcal{C} \\ (4.151)}} \frac{1}{\text{area}(\mathcal{C})} \left| \int_{\mathcal{C} \cap (P_i - H_\gamma(N))} R_j(\mathbf{x}) \, d\mathbf{x} \right| &\leq \sum_{r \geq 0} 10 \left(1 + \sqrt{\frac{\gamma}{c_1}}\right) 10 \sqrt{\frac{\gamma}{c_1}} 2^r \frac{10\gamma}{c_1 4^r} \\ \text{Case 2} & \\ &= 1000 \left( \left(\frac{\gamma}{c_1}\right)^{3/2} + \left(\frac{\gamma}{c_1}\right)^2 \right) \sum_{r \geq 0} 2^{-r} = 2000 \left( \left(\frac{\gamma}{c_1}\right)^{3/2} + \left(\frac{\gamma}{c_1}\right)^2 \right), \end{aligned} \tag{4.183}$$

a perfect analog of (4.152). This completes Case 16.

*Case 17.* The lower arc of  $P_i - H_\gamma(N)$  intersects  $\mathcal{C}$ , and the slope is less than the slope of the dominant needle  $P_{i_0} - H_\gamma(N)$ ; see Fig. 4.9.

Let  $P_{i_0} = (a_{i_0}, b_{i_0})$  and  $P_i = (a_i, b_i)$  denote the coordinates of the two points in question. By the hypothesis of Case 17, we have  $a_i > a_{i_0}$ . Write

$$h = h_i = a_i - a_{i_0} > 0 \quad \text{and} \quad v = v_i = b_i - b_{i_0},$$

where again  $h$  denotes horizontal and  $v$  denotes vertical. It is obvious from the geometry of Case 17 that  $v > 0$ . The rectangle property guarantees that  $h v \geq c_1 > 0$ .

Let  $(A_1, A_2)$  denote the coordinates of the lower left vertex of the  $j$ -cell  $\mathcal{C}$ . The intersection of the line  $x = A_1$  with the upper arc of  $P_{i_0} - H_\gamma(N)$  and the lower arc of  $P_i - H_\gamma(N)$  give two points, and the hypothesis of Case 3 implies that these intersection points are close to each other. More precisely, similar to Case 1, with  $x = 1 + a_{i_0} - A_1$ , we have the upper bound

$$\left| \left(b_{i_0} + \frac{\gamma}{x}\right) - \left(b_i - \frac{\gamma}{x + h}\right) \right| < 2 \cdot 2^{-j} \eta_2. \tag{4.184}$$



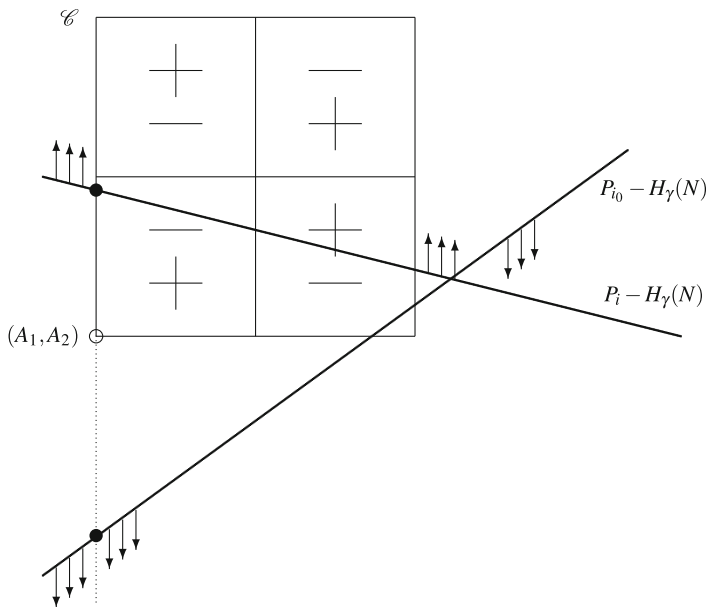


Fig. 4.9 Lower arc of  $P_i - H_\gamma(N)$  intersects  $\mathcal{C}$ , with slope less than slope of  $P_{i_0} - H_\gamma(N)$

Since  $b_i - b_{i_0} = v$ , we can rewrite (4.184) in the form

$$\left| \left( \frac{\gamma}{x} + \frac{\gamma}{x+h} \right) - v \right| < 2^{-j+1} \eta_2. \tag{4.185}$$

Note that (4.185) is an analog of (4.128) in Case 15, the only difference being that a minus sign is replaced by plus sign. This means that we can basically repeat the argument in Case 15. In fact, the plus sign helps and makes Case 17 simpler than Case 15. On the other hand, we know that the slope of the upper arc of  $\mathcal{C} \cap (P_{i_0} - H_\gamma(N))$  satisfies the inequality (4.129).

Again, if  $\eta_1$ , and so also  $\eta_2$ , is a small constant, then the upper arc of  $P_{i_0} - H_\gamma(N)$  intersects a large number of  $j$ -cells different from  $\mathcal{C}$  such that the slope is still almost equal to  $4^{-j}$ . Indeed, the horizontal size of  $\mathcal{C}$  is  $2^j \eta_1$  and, assuming that (4.129) holds, the inequality (4.130) has constant times  $1/\eta_1$  consecutive integer solutions in  $\ell$ .

Returning to (4.129) and (4.185), and then substituting  $x$  by  $x + \ell 2^j \eta_1$ , we have the respective inequalities (4.130) and

$$\left| \frac{\gamma}{x + \ell 2^j \eta_1} + \frac{\gamma}{x + \ell 2^j \eta_1 + h} - v \right| < 2^{-j+1} \eta_2. \tag{4.186}$$

If (4.129) holds, then there are at least  $\sqrt{\gamma}/10\eta_1$  consecutive integer solutions  $\ell$  of (4.130).

The basic idea is the same as in Case 15. If  $\ell$  runs through these integer solutions of (4.130) while  $\gamma, x, h$  and  $v$  remain fixed, then the function

$$\frac{\gamma}{x + \ell 2^j \eta_1} + \frac{\gamma}{x + \ell 2^j \eta_1 + h}, \tag{4.187}$$

as a function of  $\ell$ , has substantially different values, and we expect only very few of them to be very close to a fixed  $v$  in the quantitative sense of (4.186). Of course, here we assume that  $\eta_2$  is small.

Next we work out the details of this intuition. We begin by noting that (4.130) implies (4.133). Using this in (4.187), we have the good approximation

$$\frac{\gamma}{x + \ell 2^j \eta_1} + \frac{\gamma}{x + \ell 2^j \eta_1 + h} \approx \frac{\gamma}{\sqrt{\gamma} 2^j} + \frac{\gamma}{\sqrt{\gamma} 2^j + h}. \tag{4.188}$$

We now distinguish two cases. First assume that  $0 < h \leq c_1 2^{j-2} / \sqrt{\gamma}$ , where  $c_1 > 0$  is the positive constant in the rectangle property. Then the rectangle property yields

$$|v| \geq \frac{c_1}{h} \geq \frac{c_1}{\sqrt{c_1} 2^{j-1}} = 2\sqrt{c_1} 2^{-j}. \tag{4.189}$$

On the other hand, assuming that

$$\eta_2 < \frac{\sqrt{\gamma}}{2}, \tag{4.190}$$

it then follows from (4.186) and (4.188) that

$$v \leq \frac{2\gamma}{\sqrt{\gamma} 2^j} + 2^{-j+1} \eta_2 < 4\sqrt{\gamma} 2^{-j}. \tag{4.191}$$

Since (4.189) and (4.191) contradict each other, we can therefore assume that

$$h > \frac{c_1 2^{j-2}}{\sqrt{\gamma}}, \tag{4.192}$$

which is an analog of (4.138) in Case 1. Now we go back to the basic idea. We claim that if we switch  $\ell$  to  $\ell + 1$  in the function (4.187), then its value changes by at least as much as

$$\eta_1 2^{-j-2}, \tag{4.193}$$

an analog of (4.139). Indeed, (4.193) follows immediately from the routine estimate

$$\frac{1}{\sqrt{\gamma}2^j} - \frac{1}{\sqrt{\gamma}2^j + 2^j \eta_1} = \frac{1}{\sqrt{\gamma}2^j} \left( 1 - \frac{1}{1 + \eta_1/\sqrt{\gamma}} \right) \approx \frac{\eta_1}{\gamma 2^j}.$$

Let us return to (4.187) and (4.193), and apply them in (4.186). We deduce that

$$\text{at most 10 values of } \ell \text{ will satisfy both (4.130) and (4.186).} \tag{4.194}$$

As in Cases 15–16, the next step is

*A Combination of the Rectangle Property and the Pigeonhole Principle.* We recall (4.192), that  $h > c_1 2^{j-2}/\sqrt{\gamma}$ . Consider the power-of-two type decomposition

$$2^{r-1} \frac{c_1 2^{j-1}}{\sqrt{\gamma}} < h \leq 2^r \frac{c_1 2^{j-1}}{\sqrt{\gamma}}, \quad r = 0, 1, 2, \dots \tag{4.195}$$

We claim that for a fixed point  $P_{i_0} = (a_{i_0}, b_{i_0}) \in \mathcal{P}$  and for a fixed integer  $r \geq 0$ , there are at most

$$10 \cdot 2^r \tag{4.196}$$

other points  $P_i = (a_i, b_i) \in \mathcal{P}$ , with  $P_i \neq P_{i_0}$ , such that  $h = h_i = a_i - a_{i_0} > 0$  and  $v = v_i = b_i - b_{i_0}$  satisfy (4.186), thus implicitly (4.130) also, and (4.195).

To establish the bound (4.196), first note that if  $h = h_i$  satisfies (4.195), then by (4.188) and (4.195), we have

$$\frac{\gamma}{x + \ell 2^j \eta_1} + \frac{\gamma}{x + \ell 2^j \eta_1 + h} \approx \frac{\gamma}{\sqrt{\gamma} 2^j} + \frac{\gamma}{\sqrt{\gamma} 2^j + h} \approx \sqrt{\gamma} 2^{-j} \left( 1 + \frac{1}{1 + c_1 2^{r-1}/\gamma} \right),$$

so that a solution of (4.186) gives the approximation

$$v = v_i \approx \sqrt{\gamma} 2^{-j} \left( 1 + \frac{1}{1 + c_1 2^{r-1}/\gamma} \right) \pm 2^{-j+1} \eta_2. \tag{4.197}$$

Assuming

$$\eta_2 < \frac{\sqrt{\gamma}}{100}, \tag{4.198}$$

then (4.197) yields the good approximation

$$v = v_i \approx \sqrt{\gamma} 2^{-j} \left( 1 + \frac{1}{1 + c_1 2^{r-1}/\gamma} \right). \tag{4.199}$$

Suppose, contrary to the bound (4.196), that there are more than  $10 \cdot 2^r$  other points  $P_i = (a_i, b_i) \in \mathcal{P}$ , with  $P_i \neq P_{i_0}$ , such that  $h = h_i = a_i - a_{i_0} > 0$  and  $v = v_i = b_i - b_{i_0}$  satisfy (4.186), thus implicitly (4.130) also, and (4.195). Then by the Pigeonhole Principle and (4.199), there must exist two points  $P_{i_1}, P_{i_2} \in \mathcal{P}$ , with  $i_1 \neq i_2$ , such that

$$v_{i_1} \approx \sqrt{\gamma} 2^{-j} \left( 1 + \frac{1}{1 + c_1 2^{r-1}/\gamma} \right) \approx v_{i_2} \quad \text{and} \quad |h_{i_1} - h_{i_2}| \leq \frac{2^r c_1 2^{j-1}/\sqrt{\gamma}}{10 \cdot 2^r} = \frac{c_1 2^j}{20\sqrt{\gamma}}.$$

Since the product

$$\sqrt{\gamma} 2^{-j} \left( 1 + \frac{1}{1 + c_1 2^{r-1}/\gamma} \right) \cdot \frac{c_1 2^j}{2\sqrt{\gamma}} < c_1,$$

we conclude that there exists an axes-parallel rectangle of area less than  $c_1$  and which contains at least two points of  $\mathcal{P}$ , namely  $P_{i_1}$  and  $P_{i_2}$ . This contradicts the rectangle property, and establishes the bound (4.196).

If  $h = h_i$  falls into the interval (4.195), then

$$\text{slope}(\mathcal{C} \cap (P_i - H_\gamma(N))) = \frac{\gamma}{(x+h)^2} \leq \frac{\gamma}{h^2} \leq \frac{(\gamma/c_1)^2}{4^{r-2}} \cdot 4^{-j}, \tag{4.200}$$

where  $4^{-j}$  almost equals the slope of the diagonals of the  $j$ -cell  $\mathcal{C}$ . By (4.200), we have

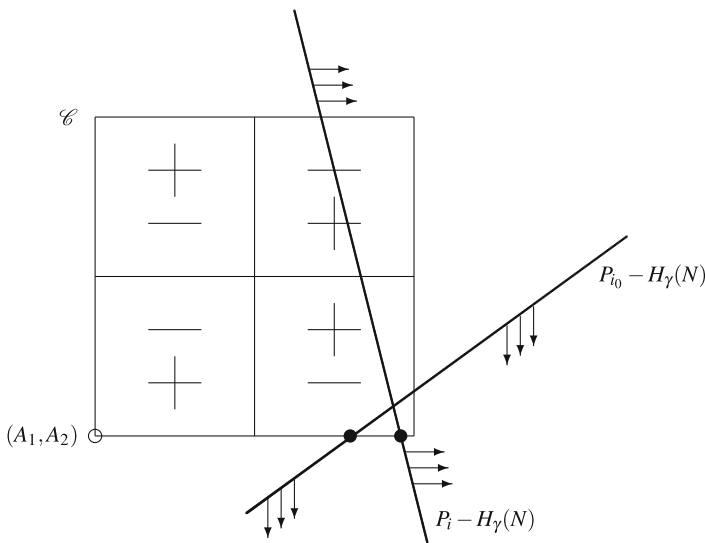
$$\frac{1}{\text{area}(\mathcal{C})} \left| \int_{\mathcal{C} \cap (P_i - H_\gamma(N))} R_j(\mathbf{x}) \, d\mathbf{x} \right| \leq \frac{10(\gamma/c_1)^2}{4^{r-2}}. \tag{4.201}$$

Furthermore, (4.201) holds for all  $j$ -cells  $\mathcal{C}$  satisfying (4.151). Let us return now to (4.126). Combining (4.194)–(4.196) and (4.201), we have

$$\begin{aligned} \sum_{\substack{P_i \in \mathcal{P} \\ i \neq i_0}} \sum_{\substack{\mathcal{C} \\ \text{(4.151)}}} \frac{1}{\text{area}(\mathcal{C})} \left| \int_{\mathcal{C} \cap (P_i - H_\gamma(N))} R_j(\mathbf{x}) \, d\mathbf{x} \right| &\leq \sum_{r \geq 0} 10 \cdot 10 \cdot 2^r \cdot \frac{10(\gamma/c_1)^2}{4^{r-2}} \\ \text{Case 3} \\ &= 16000 \left( \frac{\gamma}{c_1} \right)^2 \sum_{r \geq 0} 2^{-r} = 32000 \left( \frac{\gamma}{c_1} \right)^2. \end{aligned} \tag{4.202}$$

This completes Case 17.

*Case 18.* The lower arc of  $P_i - H_\gamma(N)$  intersects  $\mathcal{C}$ , and the slope is greater than the slope of the dominant needle  $P_{i_0} - H_\gamma(N)$ ; see Fig. 4.10.



**Fig. 4.10** Lower arc of  $P_i - H_\gamma(N)$  intersects  $\mathcal{C}$ , with slope greater than slope of  $P_{i_0} - H_\gamma(N)$

Let  $P_{i_0} = (a_{i_0}, b_{i_0})$  and  $P_i = (a_i, b_i)$  denote the coordinates of the two points in question. By the hypothesis of Case 4, we have  $a_{i_0} > a_i$ . We want positive real numbers, and write

$$h = h_i = a_{i_0} - a_i > 0 \quad \text{and} \quad v = v_i = b_i - b_{i_0} > 0,$$

where again  $h$  denotes horizontal and  $v$  denotes vertical. The rectangle property guarantees that  $hv \geq c_1 > 0$ .

Let  $(A_1, A_2)$  denote the coordinates of the lower left vertex of the  $j$ -cell  $\mathcal{C}$ . We have  $b_i > A_2 > b_{i_0}$  and  $b_i - A_2 > A_2 - b_{i_0}$ . The intersection of the line  $y = A_2$  with the upper arc of  $P_{i_0} - H_\gamma(N)$  and the lower arc of  $P_i - H_\gamma(N)$  give two points, and the hypothesis of Case 4 implies that these intersection points are relatively close to each other in the following quantitative sense. Write  $y = A_2 - b_{i_0} > 0$ . Then  $b_i - A_2 = (b_i - b_{i_0}) - y = v - y > y$ , and we have the upper bound

$$\left| \left( a_i - \frac{\gamma}{v - y} \right) - \left( a_{i_0} - \frac{\gamma}{y} \right) \right| < 2 \cdot 2^j \eta_1. \tag{4.203}$$

Since  $a_{i_0} - a_i = h > 0$ , we can rewrite (4.203) in the form

$$\left| \left( \frac{\gamma}{y} - \frac{\gamma}{v - y} \right) - h \right| < 2^{j+1} \eta_1. \tag{4.204}$$

Now we basically repeat the argument of Case 16. But, just like Case 17 is a simpler version of Case 15, Case 18 is a simpler version of Case 16. Case 18 is similar to Case 17 in the technical sense that the two critical functions

$$f_3(y) = \frac{\gamma}{y} + \frac{\gamma}{y+h} \quad \text{and} \quad f_4(y) = \frac{\gamma}{y} - \frac{\gamma}{v-y} \tag{4.205}$$

are in *synchrony*, in the sense that each is a sum of two parts that increase or decrease together as  $y$  varies.

As in Case 16, we switch the roles of the horizontal and the vertical, and focus on the reciprocal of the slope. We know that the reciprocal of the slope of the upper arc of  $\mathcal{C} \cap (P_{i_0} - H_\gamma(N))$  satisfies the inequality (4.156). We know also that if  $\eta_2$ , and so also  $\eta_1$ , is a small constant, then the upper arc of  $P_{i_0} - H_\gamma(N)$  intersects a large number of  $j$ -cells different from  $\mathcal{C}$  such that the reciprocal of the slope is still almost equal to  $4^j$ .

Returning to (4.156) and (4.204), and then substituting  $y$  by  $y + \ell 2^{-j} \eta_2$ , we have the respective inequalities (4.157) and

$$\left| \frac{\gamma}{y + \ell 2^{-j} \eta_2} - \frac{\gamma}{v - (y + \ell 2^{-j} \eta_2)} - h \right| < 2^{j+1} \eta_1. \tag{4.206}$$

If (4.156) holds, then there are at least  $\sqrt{\gamma}/10\eta_2$  consecutive integer solutions  $\ell$  of (4.157).

The basic idea is the same as in Case 16. If  $\ell$  runs through these integer solutions of (4.157) while  $\gamma, x, h$  and  $v$  remain fixed, then the function

$$\frac{\gamma}{y + \ell 2^{-j} \eta_2} - \frac{\gamma}{v - (y + \ell 2^{-j} \eta_2)}, \tag{4.207}$$

as a function of  $\ell$ , has substantially different values, and we expect only very few of them to be very close to a fixed  $h$  in the quantitative sense of (4.157). Of course, here we assume that  $\eta_1$  is small.

Next we work out the details of this intuition. We begin by noting that (4.157) implies (4.160). Since the functions  $f_3(y)$  and  $f_4(y)$  given by (4.205) are in synchrony, we can basically repeat the argument of (4.187), (4.193) and (4.194) in Case 3, and conclude that if we switch  $\ell$  to  $\ell + 1$  in the function (4.207), then its value changes by at least as much as

$$\eta_2 2^{j-2},$$

an analog of (4.171) and (4.193). Thus we deduce that

$$\text{at most 10 values of } \ell \text{ will satisfy both (4.157) and (4.206).} \tag{4.208}$$

As in Cases 15–17, the next step is

*A Combination of the Rectangle Property and the Pigeonhole Principle.* In this case, since  $v - (y + \ell 2^{-j} \eta_2) > y + \ell 2^{-j} \eta_2$ , we have

$$v > 2(y + \ell 2^{-j} \eta_2). \quad (4.209)$$

In view of (4.160), we can assume that

$$v > \sqrt{\frac{6\gamma}{7}} 2^{-j+1}.$$

Consider the power-of-two type decomposition

$$2^{r-1} \sqrt{\frac{6\gamma}{7}} 2^{-j+2} < v \leq 2^r \sqrt{\frac{6\gamma}{7}} 2^{-j+2}, \quad r = 0, 1, 2, \dots \quad (4.210)$$

We claim that for a fixed point  $P_{i_0} = (a_{i_0}, b_{i_0}) \in \mathcal{P}$  and for a fixed integer  $r \geq 0$ , there are at most

$$\frac{100\gamma 2^r}{c_1} \quad (4.211)$$

other points  $P_i = (a_i, b_i) \in \mathcal{P}$ , with  $P_i \neq P_{i_0}$ , such that  $h = h_i = a_{i_0} - a_i > 0$  and  $v = v_i = b_i - b_{i_0} > 0$  satisfy (4.206), thus implicitly (4.157) also, and (4.210).

To establish the bound (4.211), first note that if  $v = v_i$  satisfies (4.210), then by (4.160), (4.206) and (4.209), and assuming that

$$\eta_1 < \frac{\sqrt{\gamma}}{4}, \quad (4.212)$$

we have

$$h = h_i < \frac{\gamma}{y + \ell 2^{-j} \eta_2} + 2^{j+1} \eta_1 \leq \frac{\gamma}{2^{-j} \sqrt{6\gamma/7}} + 2^{j+1} \eta_1 \leq 2\sqrt{\gamma} 2^j. \quad (4.213)$$

Suppose, contrary to the bound (4.211), that there are more than  $100\gamma 2^r / c_1$  other points  $P_i = (a_i, b_i) \in \mathcal{P}$ , with  $P_i \neq P_{i_0}$ , such that  $h = h_i = a_{i_0} - a_i > 0$  and  $v = v_i = b_i - b_{i_0} > 0$  satisfy (4.206), thus implicitly (4.157) also, and (4.210). Then by the Pigeonhole Principle and (4.213), there must exist two points  $P_{i_1}, P_{i_2} \in \mathcal{P}$ , with  $i_1 \neq i_2$ , such that

$$\max\{h_{i_1}, h_{i_2}\} \leq 2\sqrt{\gamma} 2^j \quad \text{and} \quad |v_{i_1} - v_{i_2}| \leq \frac{2^r \sqrt{6\gamma/7} 2^{-j+2}}{100\gamma 2^r / c_1} = \frac{c_1 \sqrt{6/7}}{25\sqrt{\gamma}} 2^{-j}.$$

Since the product

$$\sqrt{\gamma}2^j \cdot \frac{c_1\sqrt{6/7}}{\sqrt{\gamma}}2^{-j} = \sqrt{\frac{6}{7}}c_1 < c_1,$$

we conclude that there exists an axes-parallel rectangle of area less than  $c_1$  and which contains at least two points of  $\mathcal{P}$ , namely  $P_{i_1}$  and  $P_{i_2}$ . This contradicts the rectangle property, and establishes the bound (4.211).

If  $v = v_i$  falls into the interval (4.210), then

$$\frac{1}{\text{slope}(\mathcal{C} \cap (P_i - H_\gamma(N)))} = \frac{\gamma}{(y+v)^2} \leq \frac{\gamma}{v^2} \approx \frac{1}{4^r} \cdot 4^j, \tag{4.214}$$

where  $4^j$  almost equals the reciprocal of the slope of the diagonals of the  $j$ -cell  $\mathcal{C}$ . By (4.214), we have

$$\frac{1}{\text{area}(\mathcal{C})} \left| \int_{\mathcal{C} \cap (P_i - H_\gamma(N))} R_j(\mathbf{x}) \, d\mathbf{x} \right| \leq \frac{10}{4^r}. \tag{4.215}$$

Furthermore, (4.215) holds for all  $j$ -cells  $\mathcal{C}$  satisfying (4.151). Let us return now to (4.126). Combining (4.208), (4.210), (4.211) and (4.215) we have

$$\begin{aligned} \sum_{\substack{P_i \in \mathcal{P} \\ i \neq i_0 \\ \text{Case 4}}} \sum_{\mathcal{C}} \frac{1}{\text{area}(\mathcal{C})} \left| \int_{\mathcal{C} \cap (P_i - H_\gamma(N))} R_j(\mathbf{x}) \, d\mathbf{x} \right| &\leq \sum_{r \geq 0} 10 \cdot 100 \frac{\gamma}{c_1} 2^r \cdot \frac{10}{4^r} \\ &= 10000 \frac{\gamma}{c_1} \sum_{r \geq 0} 2^{-r} = 20000 \frac{\gamma}{c_1}. \end{aligned} \tag{4.216}$$

This completes Case 18.

### 4.8 Completing the Proof of Theorem 12

In this section, we shall finally complete the proof of Proposition 13. Let us return to (4.125) and (4.126). We are now ready to clarify the technical details of the single term domination.

Let  $P_{i_0} \in \mathcal{P}$  and  $j \in \mathcal{J}$  be arbitrary.

*Property 19.* The slope  $\gamma/x^2$  of the hyperbolic needle  $P_{i_0} - H_\gamma(N)$  satisfies

$$\frac{5}{6}4^{-j} \leq \frac{\gamma}{x^2} \leq \frac{7}{6}4^{-j}. \tag{4.217}$$



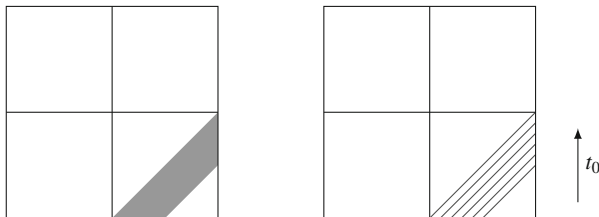


Fig. 4.11 Short vertical translations

Note that (4.217) holds if and only if

$$\sqrt{\frac{6\gamma}{7}}2^j \leq x \leq \sqrt{\frac{6\gamma}{5}}2^j,$$

and this is an interval of length greater than  $\sqrt{\gamma}2^j/6$ . Since a  $j$ -cell  $\mathcal{C}$  has horizontal side  $\eta_1 2^j$ , there are more than

$$\frac{\sqrt{\gamma}2^j/6}{\eta_1 2^j} = \frac{\sqrt{\gamma}}{6\eta_1}$$

$j$ -cells  $\mathcal{C}$  with the slope of the intersection  $\mathcal{C} \cap (P_{i_0} - H_\gamma(N))$  satisfying Property 19.

It would be not too difficult to prove directly, by using some familiar arguments from uniform distribution, that among these more than  $\sqrt{\gamma}/6\eta_1$   $j$ -cells  $\mathcal{C}$ , at least 1 % has the following additional property.

*Property 20.* The hyperbolic needle  $P_{i_0} - H_\gamma(N)$  intersects only the lower right subrectangle of  $\mathcal{C}$ , and the intersection is a large triangle, meaning that the area is at least  $\frac{1}{32}$  the area of  $\mathcal{C}$ , i.e. the area is at least  $\eta_1 \eta_2/32$ .

It is technically simpler, however, to force Property 20 in an indirect way, by using the trick of short vertical translations; see Fig. 4.11. This geometric trick was already mentioned at the end of Sect. 4.5.

More precisely, for every real number  $t_0$  satisfying  $0 < t_0 < 1$ , consider all  $j$ -cells  $\mathcal{C}$  such that, with  $B = [0, M - N] \times [\gamma, M - \gamma]$ , we have

$$\mathcal{C} \cap (P_{i_0} + (0, t_0) - H_\gamma(N)) \subset B \tag{4.218}$$

and

$$\frac{5}{6}4^{-j} \leq \text{slope}(\mathcal{C} \cap (P_{i_0} + (0, t_0) - H_\gamma(N))) \leq \frac{7}{6}4^{-j}. \tag{4.219}$$

Simple geometric consideration shows that for, say, at least 5 % of the pairs  $(t_0, \mathcal{C})$ , where  $\mathcal{C}$  satisfies (4.218) and (4.219),  $\mathcal{C} \cap (P_{i_0} + (0, t_0) - H_\gamma(N))$  also satisfies

Property 20, i.e.  $P_{i_0} + (0, t_0) - H_\gamma(N)$  intersects only the lower right subrectangle of  $\mathcal{C}$ , and the intersection is a large triangle of area at least  $\eta_1 \eta_2 / 32$ .

For the proof of the positive direction (4.89), we choose the pattern  $+-$  in every  $j$ -cell  $\mathcal{C}$  satisfying (4.218) and (4.219). Naturally, we choose the opposite pattern  $-+$  for the negative direction (4.90). Then

$$\int_{\mathcal{C} \cap (P_{i_0} + (0, t_0) - H_\gamma(N))} R_j(\mathbf{x}) \, d\mathbf{x} \geq \frac{\eta_1 \eta_2}{32}. \tag{4.220}$$

Finally, if the  $j$ -cell  $\mathcal{C}$  does not satisfy both (4.218) and (4.219), then we choose the pattern 0. Therefore, by (4.220) and summarizing Cases 1–4, we have

$$\begin{aligned} & \int_0^1 \left( \sum_{j \in \mathcal{J}} \sum_{P_{i_0} \in \mathcal{P}} \int_{P_{i_0} + (0, t_0) - H_\gamma(N)} R_j(\mathbf{x}) \, d\mathbf{x} \right) dt_0 \\ & \geq \sum_{j \in \mathcal{J}} \sum_{P_{i_0} \in \mathcal{P}} \left( \frac{1}{20} \cdot \frac{\sqrt{\gamma}}{6\eta_1} \cdot \frac{\eta_1 \eta_2}{32} \right. \\ & \quad \left. - \sum_{\substack{P_i \in \mathcal{P} \\ i \neq i_0}} \sum_{\mathcal{C}} \int_0^1 \left| \int_{\mathcal{C} \cap (P_i + (0, t_0) - H_\gamma(N))} R_j(\mathbf{x}) \, d\mathbf{x} \right| dt_0 \right) \\ & \geq \sum_{j \in \mathcal{J}} \sum_{P_{i_0} \in \mathcal{P}} \left( \frac{\sqrt{\gamma} \eta_2}{3840} - \eta_1 \eta_2 \left( 4000 \left( \left( \frac{\gamma}{c_1} \right)^{3/2} + \left( \frac{\gamma}{c_1} \right)^2 \right) \right. \right. \\ & \quad \left. \left. + 32000 \left( \frac{\gamma}{c_1} \right)^2 + 20000 \frac{\gamma}{c_1} \right) \right), \tag{4.221} \end{aligned}$$

where the summation over  $P_{i_0} \in \mathcal{P}$  is under the restriction

$$P_{i_0} + (0, t_0) - H_\gamma(N) \subset B \quad \text{for all } t_0 \text{ satisfying } 0 < t_0 < 1, \tag{4.222}$$

the summation over  $\mathcal{C}$  is under the restriction (4.219), the summation over  $P_i \in \mathcal{P}$  with  $i \neq i_0$  encompass Cases 15–18, and finally the factor  $\frac{1}{20}$  comes from the 5% mentioned earlier. Furthermore, we have used in the last step the inequalities (4.152), (4.183), (4.202) and (4.216) for every  $t_0$  satisfying  $0 < t_0 < 1$ .

In our discussion in Sects. 4.6 and 4.7, we have made some assumptions on  $\eta_1$  and  $\eta_2$ . Corresponding to Cases 15–18, we have assumed respectively that

$$\begin{aligned} \eta_2 &< \min \left\{ \frac{\sqrt{c_1}}{2}, \frac{1}{8(1/\sqrt{\gamma} + 1/\sqrt{c_1})} \right\}, \\ \eta_1 &< \min \left\{ \frac{c_1}{2\sqrt{\gamma}}, \frac{\sqrt{\gamma}}{8}, \frac{\sqrt{c_1}}{2}, \frac{1}{8(1/\sqrt{\gamma} + 1/\sqrt{c_1})} \right\}, \\ \eta_2 &< \min \left\{ \frac{\sqrt{\gamma}}{2}, \frac{\sqrt{\gamma}}{100} \right\}, \\ \eta_1 &< \frac{\sqrt{\gamma}}{4}; \end{aligned}$$

see (4.137), (4.147), (4.164), (4.166), (4.169), (4.179), (4.190), (4.198) and (4.212).  
 Since

$$\frac{1}{1/\sqrt{\gamma} + 1/\sqrt{c_1}} \geq \frac{\sqrt{\gamma} + \sqrt{c_1}}{2},$$

we can guarantee all of the above requirements on  $\eta_1$  and  $\eta_2$  by imposing the single inequality

$$\max\{\eta_1, \eta_2\} < \min \left\{ \frac{\sqrt{\gamma}}{100}, \frac{\sqrt{c_1}}{8}, \frac{c_1}{2\sqrt{\gamma}} \right\}. \tag{4.223}$$

Let us return to (4.221). We have

$$\begin{aligned} &\frac{\sqrt{\gamma}\eta_2}{3840} - \eta_1\eta_2 \left( 4000 \left( \left( \frac{\gamma}{c_1} \right)^{3/2} + \left( \frac{\gamma}{c_1} \right)^2 \right) + 32000 \left( \frac{\gamma}{c_1} \right)^2 + 20000 \frac{\gamma}{c_1} \right) \\ &\geq \frac{\sqrt{\gamma}\eta_2}{7680}, \end{aligned} \tag{4.224}$$

assuming that (4.223) holds and  $\eta_1$  satisfies the additional inequality

$$\frac{1}{\eta_1} \geq \frac{10^8}{\sqrt{\gamma}} \left( \left( \frac{\gamma}{c_1} \right) + \left( \frac{\gamma}{c_1} \right)^2 \right). \tag{4.225}$$

Since  $\eta_1$  and  $\eta_2$  are almost equal, in view of (4.100), we can satisfy both (4.223) and (4.225) by the choice

$$\eta_1 \approx \eta_2 = \min \left\{ \frac{\sqrt{\gamma}}{200}, \frac{\sqrt{c_1}}{10}, \frac{10^{-8}c_1}{2\sqrt{\gamma}}, \frac{10^{-8}c_1^2}{2\gamma^{3/2}} \right\}. \tag{4.226}$$

Substituting (4.226) in (4.224) and then returning to (4.221), we have

$$\int_0^1 \left( \sum_{j \in \mathcal{J}} \sum_{P_{i_0} \in \mathcal{P}} \int_{P_{i_0} + (0, t_0) - H_\gamma(N)} R_j(\mathbf{x}) \, d\mathbf{x} \right) dt_0 \geq \sum_{j \in \mathcal{J}} \sum_{P_{i_0} \in \mathcal{P}} \frac{\sqrt{\gamma} \eta_2}{7680}, \tag{4.222}$$

where  $i_0$  is now a dummy variable. Clearly there exists  $t_0$ , satisfying  $0 < t_0 < 1$ , such that

$$\sum_{j \in \mathcal{J}} \sum_{P_i \in \mathcal{P}} \int_{P_i + (0, t_0) - H_\gamma(N)} R_j(\mathbf{x}) \, d\mathbf{x} \geq \sum_{j \in \mathcal{J}} \sum_{P_i \in \mathcal{P}} \frac{\sqrt{\gamma} \eta_2}{7680}. \tag{4.227}$$

(4.228)

Note that in (4.227), we have substituted the dummy variable  $i_0$  by  $i$ , together with a corresponding summation restriction

$$P_i + (0, t_0) - H_\gamma(N) \subset B \quad \text{for all } t_0 \text{ satisfying } 0 < t_0 < 1. \tag{4.228}$$

Next we return to (4.118), and replace the point set  $\mathcal{P}$  by the translated point set  $\mathcal{P} + (0, t_0)$ . Then Lemma 14 gives

$$\begin{aligned} & \frac{1}{(M - N)(M - 2\gamma)} \int_0^{M-N} \int_\gamma^{M-\gamma} \Delta(\mathbf{x}) T(\mathbf{x}) \, d\mathbf{x} \\ &= \sum_{P_i \in \mathcal{P}} \frac{\text{area}(B \cap (P_i + (0, t_0) - H_\gamma(N)))}{(M - N)(M - 2\gamma)} - \delta \cdot \text{area}(H_\gamma(N)) \\ &+ \rho \sum_{j \in \mathcal{J}} \sum_{P_i \in \mathcal{P}} \frac{1}{(M - N)(M - 2\gamma)} \int_{P_i + (0, t_0) - H_\gamma(N)} R_j(\mathbf{x}) \, d\mathbf{x} + E_1, \end{aligned} \tag{4.229}$$

where the error  $E_1$  satisfies

$$|E_1| \leq \frac{|\mathcal{P}|}{(M - N)(M - 2\gamma)} \cdot 8\sqrt{\gamma}(\eta_1 + \eta_2) \cdot \frac{(n + 1)\rho^2}{\sqrt{2} - 1 - \rho}. \tag{4.230}$$

Recall that  $\mathcal{P}$  is a finite subset of the square  $[0, M]^2$  with cardinality  $|\mathcal{P}| = \delta M^2$ . Since  $0 < t_0 < 1$ , the rectangle property implies, via elementary calculations, that the condition

$$P_i + (0, t_0) - H_\gamma(N) \subset B = [0, M - N] \times [\gamma, M - \gamma] \tag{4.231}$$

holds for all but at most

$$\frac{(2N + 4\gamma + 1)M}{c_1} \tag{4.232}$$

points  $P_i \in \mathcal{P}$ . Thus

$$\begin{aligned}
& \sum_{P_i \in \mathcal{P}} \frac{\text{area}(B \cap (P_i + (0, t_0) - H_\gamma(N)))}{(M - N)(M - 2\gamma)} - \delta \cdot \text{area}(H_\gamma(N)) \\
&= \frac{\delta M^2 + \theta c_1^{-1}(2N + 4\gamma + 1)M}{(M - N)(M - 2\gamma)} \cdot \text{area}(H_\gamma(N)) - \delta \cdot \text{area}(H_\gamma(N)) \\
&= \left( \frac{M^2}{(M - N)(M - 2\gamma)} - 1 \right) \delta \cdot \text{area}(H_\gamma(N)) \\
&\quad + \theta \frac{c_1^{-1}(2N + 4\gamma + 1)M}{(M - N)(M - 2\gamma)} \cdot \text{area}(H_\gamma(N)),
\end{aligned}$$

with some constant  $\theta$  satisfying  $-1 \leq \theta \leq 1$ . Since  $\text{area}(H_\gamma(N)) = 2\gamma \log N$ , it then follows that

$$\begin{aligned}
& \left| \sum_{P_i \in \mathcal{P}} \frac{\text{area}(B \cap (P_i + (0, t_0) - H_\gamma(N)))}{(M - N)(M - 2\gamma)} - \delta \cdot \text{area}(H_\gamma(N)) \right| \\
&\leq \frac{3N + 6\gamma + 1}{(M - N)(M - 2\gamma)} \cdot 2\gamma \log N.
\end{aligned} \tag{4.233}$$

Combining (4.229), (4.230) and (4.233), we deduce that

$$\begin{aligned}
& \frac{1}{(M - N)(M - 2\gamma)} \int_0^{M-N} \int_\gamma^{M-\gamma} \Delta(\mathbf{x}) T(\mathbf{x}) \, \mathbf{d}\mathbf{x} \\
&= \rho \sum_{j \in \mathcal{J}} \sum_{P_i \in \mathcal{P}} \frac{1}{(M - N)(M - 2\gamma)} \int_{P_i + (0, t_0) - H_\gamma(N)} R_j(\mathbf{x}) \, \mathbf{d}\mathbf{x} + E_2,
\end{aligned} \tag{4.234}$$

where the error  $E_2$  satisfies

$$\begin{aligned}
|E_2| &\leq \frac{|\mathcal{P}|}{(M - N)(M - 2\gamma)} \cdot 8\sqrt{\gamma}(\eta_1 + \eta_2) \cdot \frac{(n + 1)\rho^2}{\sqrt{2} - 1 - \rho} \\
&\quad + \frac{3N + 6\gamma + 1}{(M - N)(M - 2\gamma)} \cdot 2\gamma \log N.
\end{aligned} \tag{4.235}$$

Combining (4.227), (4.234) and (4.235), we then conclude that

$$\begin{aligned}
& \frac{1}{(M - N)(M - 2\gamma)} \int_0^{M-N} \int_\gamma^{M-\gamma} \Delta(\mathbf{x}) T(\mathbf{x}) \, \mathbf{d}\mathbf{x} \\
&\geq \rho \sum_{j \in \mathcal{J}} \frac{1}{(M - N)(M - 2\gamma)} \sum_{P_i \in \mathcal{P}} \frac{\sqrt{\gamma}\eta_2}{7680}
\end{aligned} \tag{4.228}$$

$$\begin{aligned} & -\frac{|\mathcal{P}|}{(M-N)(M-2\gamma)} \cdot 8\sqrt{\gamma}(\eta_1 + \eta_2) \cdot \frac{(n+1)\rho^2}{\sqrt{2}-1-\rho} \\ & -\frac{3N+6\gamma+1}{(M-N)(M-2\gamma)} \cdot 2\gamma \log N. \end{aligned} \tag{4.236}$$

Recall that  $\mathcal{J}$  is an interval of integers satisfying (4.114), so that

$$|\mathcal{J}| \geq (n+1) - \log_2 \left( \gamma + \frac{1}{\gamma} \right).$$

On the other hand, it follows from (4.231) and (4.232) that

$$\sum_{\substack{P_i \in \mathcal{P} \\ (4.228)}} 1 \geq \delta M^2 - \frac{(2N+4\gamma+1)M}{c_1}.$$

Thus

$$\begin{aligned} & \sum_{j \in \mathcal{J}} \frac{1}{(M-N)(M-2\gamma)} \sum_{\substack{P_i \in \mathcal{P} \\ (4.228)}} 1 \\ & \geq \left( (n+1) - \log_2 \left( \gamma + \frac{1}{\gamma} \right) \right) \left( \delta - \frac{2N+4\gamma+1}{c_1 M} \right). \end{aligned} \tag{4.237}$$

Let us now return to (4.236). If  $\rho$  is small, then  $\rho^2$  is negligible compared to  $\rho$ . Let  $\rho = 10^{-6}$ , say. Substituting this and the estimate (4.237) into (4.236), and assuming that  $N$  and  $M/N$  are both large, we deduce that

$$\begin{aligned} & \frac{1}{\text{area}(B)} \int_B \Delta(\mathbf{x})T(\mathbf{x}) \, d\mathbf{x} \\ & \geq \rho \left( (n+1) - \log_2 \left( \gamma + \frac{1}{\gamma} \right) \right) \left( \delta - \frac{2N+4\gamma+1}{c_1 M} \right) \frac{\sqrt{\gamma}\eta_2}{10^4}. \end{aligned}$$

More precisely, the assumptions on  $N$  and  $M$  are given by (4.84) and (4.85), and the choice for  $n$  is made precise by

$$\frac{N}{2} < 2^n \leq N.$$

These choices, together with the definition (4.226) for  $\eta_2$ , ensure that

$$\frac{1}{\text{area}(B)} \int_B \Delta(\mathbf{x})T(\mathbf{x}) \, d\mathbf{x} \geq \delta' \log N, \tag{4.238}$$

where  $\delta' = \delta'(c_1, \gamma, \delta) > 0$  is a positive constant independent of  $N$  and  $M$ , and defined by (4.83) and (4.84).

It now follows from (4.238) that there exists a translated copy  $\mathbf{x}_1 + H_\gamma(N)$  of the hyperbolic needle  $H_\gamma(N)$  such that  $\mathbf{x}_1 + H_\gamma(N) \subset [0, M]^2$  and

$$|\mathcal{P} \cap (\mathbf{x}_1 + H_\gamma(N))| \geq 2\delta\gamma \log N + \delta' \log N.$$

This establishes the inequality (4.89). The proof of the other inequality (4.90) is the same, except that we replace the pattern  $+-$  by the opposite pattern  $-+$ .

Thus the long proof of Proposition 13 is complete. This also completes the proof of Theorem 12.

### 4.9 Yet Another Generalization of Theorem 3

Let  $\alpha > 0$ ,  $0 \leq \beta < 1$  and  $\gamma > 0$  be arbitrary but fixed real numbers, and let  $f(\alpha; \beta; \gamma; N)$  denote the number of integral solutions of the diophantine inequality<sup>14</sup>

$$\|n\alpha - \beta\| < \frac{\gamma}{n}, \quad 1 \leq n \leq N.$$

This inequality motivates the hyperbolic region

$$|y - \beta| < \frac{\gamma}{x}, \quad 1 \leq x \leq N,$$

which has area  $2\gamma \log N$ .

Let us return to the special case  $\alpha = \sqrt{2}$ . Combining Lemmas 1 and 2, we have

$$\int_0^1 f(\sqrt{2}; \beta; \gamma; N) \, d\beta = 2\gamma \log N + O(1), \tag{4.239}$$

and for an arbitrary subinterval  $[a, b]$  with  $0 \leq a < b \leq 1$ , we have the limit formula

$$\lim_{N \rightarrow \infty} \frac{\frac{1}{b-a} \int_a^b f(\sqrt{2}; \beta; \gamma; N) \, d\beta}{\log N} = 2\gamma. \tag{4.240}$$

There is a straightforward generalization of (4.239) and (4.240) for arbitrary  $\alpha > 0$ , and the proof is the same. We have

$$\int_0^1 f(\alpha; \beta; \gamma; N) \, d\beta = 2\gamma \log N + O(1), \tag{4.241}$$

---

<sup>14</sup>Note that the special case  $\alpha = \sqrt{2}$  was introduced in Sect. 4.1; see (4.28).

and for an arbitrary subinterval  $[a, b]$  with  $0 \leq a < b \leq 1$ , we have the limit formula

$$\lim_{N \rightarrow \infty} \frac{\frac{1}{b-a} \int_a^b f(\alpha; \beta; \gamma; N) \, d\beta}{\log N} = 2\gamma. \tag{4.242}$$

The formulas (4.239)–(4.242) express the almost trivial geometric fact that the average number of lattice points contained in all the translated copies of a given region equals the area of the region; see Lemma 5. It is natural, therefore, to study the limit

$$\lim_{N \rightarrow \infty} \frac{f(\alpha; \beta; \gamma; N)}{2\gamma \log N}. \tag{4.243}$$

The case of rational  $\alpha$  in (4.243) is trivial. Indeed, if  $N \rightarrow \infty$ , then the function  $f(\alpha; \beta; \gamma; N)$  remains bounded for all but a finite number of values of  $\beta = \beta(\alpha)$  in the unit interval. When  $f(\alpha; \beta; \gamma; N)$  tends to infinity, it behaves like a linear function  $c_{27}N$ , which is much faster than the logarithmic function  $\log N$ .

If  $\alpha$  is irrational, then we have the following non-trivial result, which can be considered a far-reaching generalization of Theorem 3.

**Theorem 21.** *Let  $\alpha > 0$  be an arbitrary irrational, and let  $\gamma > 0$  be an arbitrary real number. There are continuum many divergence points  $\beta^* = \beta^*(\alpha, \gamma) \in [0, 1)$  such that*

$$\limsup_{n \rightarrow \infty} \frac{f(\alpha; \beta^*; \gamma; n)}{\log n} > \liminf_{n \rightarrow \infty} \frac{f(\alpha; \beta^*; \gamma; n)}{\log n}.$$

To prove Theorem 21, we can clearly assume that  $0 < \alpha < 1$ . We need the continued fractions

$$\alpha = \frac{1}{a_1 + \frac{1}{a_2 + \frac{1}{a_3 + \dots}}} = [a_1, a_2, a_3, \dots].$$

For irrational  $\alpha$ , the digits  $a_1, a_2, a_3, \dots$  form an infinite sequence, with  $a_i \geq 1$  for all  $i \geq 1$ . For  $k \geq 2$ , the fractions

$$\frac{p_k}{q_k} = [a_1, \dots, a_k]$$

are known as the convergents to  $\alpha$ . It is well known that  $p_k, q_k$  are generated by the recurrence relations

$$p_k = a_k p_{k-1} + p_{k-2}, \quad q_k = a_k q_{k-1} + q_{k-2}, \tag{4.244}$$

with the convention that  $p_0 = 0, q_0 = 1, p_1 = 1$  and  $q_1 = a_1$ .



Another well-known fact about the convergents is the inequality

$$\left| \alpha - \frac{p_{k-1}}{q_{k-1}} \right| \leq \frac{1}{q_{k-1}q_k},$$

which clearly implies

$$\|q_{k-1}\alpha\| < \frac{1}{a_k q_{k-1}}. \tag{4.245}$$

Write  $n = \ell q_{k-1}$ . Then by (4.245), we have

$$\|n\alpha\| = \|\ell q_{k-1}\alpha\| < \frac{\ell}{a_k q_{k-1}} = \frac{\ell^2}{a_k \ell q_{k-1}} = \frac{\ell^2}{a_k n},$$

and so  $\|n\alpha\| < \gamma/n$  holds whenever  $\ell^2/a_k \leq \gamma$ , i.e. whenever

$$1 \leq \ell \leq \sqrt{\gamma a_k}. \tag{4.246}$$

Now let

$$N_k = \lfloor \sqrt{\gamma a_k} \rfloor q_{k-1}, \tag{4.247}$$

where  $\lfloor z \rfloor$  denotes the lower integral part of a real number  $z$ . It then follows from (4.246) that the homogeneous diophantine inequality  $\|n\alpha\| < \gamma/n$  has at least

$$\sum_{i=1}^k \lfloor \sqrt{\gamma a_i} \rfloor$$

integer solutions  $n$  satisfying  $1 \leq n \leq N_k$ . Formally, we therefore have

$$f(\alpha; \beta = 0; \gamma; N_k) \geq \sum_{i=1}^k \lfloor \sqrt{\gamma a_i} \rfloor. \tag{4.248}$$

We distinguish two cases, and start with the much harder one.

*Case 22.* For all sufficiently large values of  $k$ , we have

$$\sum_{i=1}^k \lfloor \sqrt{\gamma a_i} \rfloor \leq 100 \cdot 2\gamma \log N_k. \tag{4.249}$$

We proceed in four steps.

*Step 1.* The crucial first step in the argument is to show that the condition (4.249) implies the exponential upper bound

$$\prod_{i=1}^k (a_i + 1) \leq e^{c'k} \tag{4.250}$$

for all sufficiently large values of  $k$ , where  $c' = c'(\gamma)$  is a finite constant independent of  $k$ .

To derive (4.250), we use the well-known principle that the exponential functions grow faster than polynomials, in the form of an elementary inequality as follows.

**Lemma 23.** *For any fixed positive  $c > 0$ , the inequality*

$$(x + 1)^c \leq (8c^2 e^{-2})^c e^{\sqrt{x}}$$

*holds for every  $x \geq 1$ .*

*Proof.* We start with the trivial observation that  $x + 1 \leq 2x$  for all  $x \geq 1$ , which leads us to the function  $g(x) = (2x)^c e^{-\sqrt{x}}$ , which we wish to maximize. It is easy to compute the derivative of  $g(x)$  and show that its value is maximized when  $x = 4c^2$ . The desired inequality follows from  $(x + 1)^c e^{-\sqrt{x}} \leq g(x) \leq g(4c^2)$ .  $\square$

By repeated application of (4.244), we have

$$\begin{aligned} q_{k-1} &= a_{k-1}q_{k-2} + q_{k-3} \leq (a_{k-1} + 1)q_{k-2} \\ &\leq (a_{k-1} + 1)(a_{k-2} + 1)q_{k-3} \leq \dots \leq \prod_{i=1}^{k-1} (a_i + 1). \end{aligned} \tag{4.251}$$

Combining this with (4.247) and (4.249), we have

$$\begin{aligned} \sum_{i=1}^k (\sqrt{\gamma a_i} - 1) &\leq 100 \cdot 2\gamma \left( \log \sqrt{\gamma} + \log \sqrt{a_k} + \log \prod_{i=1}^{k-1} (a_i + 1) \right) \\ &\leq 200\gamma \left( \log \sqrt{\gamma} + \log \prod_{i=1}^k (a_i + 1) \right). \end{aligned} \tag{4.252}$$

Applying the exponential function, the inequality (4.252) becomes

$$\prod_{i=1}^k e^{\sqrt{\gamma a_i} - 1} \leq \gamma^{100\gamma} \prod_{i=1}^k (a_i + 1)^{200\gamma}, \tag{4.253}$$

and this inequality holds for all sufficiently large  $k$ , i.e. for all  $k \geq k_0$ .

Applying Lemma 23 with  $c = 400\sqrt{\gamma}$  and  $x = a_i$  for each  $i = 1, 2, \dots, k + 1$ , and then multiplying these inequalities together, we obtain

$$\prod_{i=1}^k (a_i + 1)^{400\sqrt{\gamma}} \leq (800\sqrt{\gamma})^{800\sqrt{\gamma}k} \prod_{i=1}^k e^{\sqrt{a_i}}.$$

Raising this to the  $\sqrt{\gamma}$ -th power, we have

$$\prod_{i=1}^k (a_i + 1)^{400\gamma} \leq (800\gamma)^{800\gamma k} \prod_{i=1}^k e^{\sqrt{\gamma}a_i}. \tag{4.254}$$

We next combine (4.253) and (4.254) to obtain

$$\prod_{i=1}^k (a_i + 1)^{400\gamma} \leq (800\gamma)^{800\gamma k} e^k \gamma^{100\gamma} \prod_{i=1}^k (a_i + 1)^{200\gamma},$$

which, on removing a common factor and then taking  $200\gamma$ -th root, becomes

$$\prod_{i=1}^k (a_i + 1) \leq (800\gamma)^{4k} e^{k/200\gamma} \sqrt{\gamma} = \sqrt{\gamma}((800\gamma)^4 e^{1/200\gamma})^k. \tag{4.255}$$

Since this holds for all  $k \geq k_0$ , the inequality (4.250) follows.

*Step 2.* We shall next show that small digit  $a_i$  implies a local rectangle property. It follows from (4.255) that, for all sufficiently large  $k$ ,

$$a_i + 1 \leq (1000\gamma)^8 e^{\frac{1}{100\gamma}} \tag{4.256}$$

holds for at least  $k/2$  values of  $i$  in  $1 \leq i \leq k$ . In other words, at least half of the continued fraction digits  $a_i$  of  $\alpha$  are small, less than a constant depending only on  $\gamma$ , in the precise quantitative sense of (4.256).

Next we show that, for every small digit  $a_i$ , the rectangle property must hold locally, in some power-of-two range around  $q_i$ . To prove this, we basically repeat the proof of Lemma 4, and use some facts from the theory of continued fractions; see Lemma 24 below. The details go as follows.

As in the proof of Lemma 4, we consider a rectangle of slope  $1/\alpha$  and which contains two lattice points  $P = (k, \ell)$  and  $Q = (m, n)$ ; in fact, assume that  $P$  and  $Q$  are two vertices of the rectangle. We denote the vector from  $P$  to  $Q$  by  $\mathbf{v} = (m - k, n - \ell)$ , and consider the two perpendicular unit vectors

$$\mathbf{e}_1 = \left( \frac{\alpha}{\sqrt{1 + \alpha^2}}, \frac{1}{\sqrt{1 + \alpha^2}} \right) \quad \text{and} \quad \mathbf{e}_2 = \left( \frac{1}{\sqrt{1 + \alpha^2}}, -\frac{\alpha}{\sqrt{1 + \alpha^2}} \right).$$

Then the two sides  $a$  and  $b$  of the rectangle can be expressed in terms of the inner products  $\mathbf{e}_1 \cdot \mathbf{v}$  and  $\mathbf{e}_2 \cdot \mathbf{v}$ . We have

$$a = |\mathbf{e}_1 \cdot \mathbf{v}| = \frac{|p\alpha + q|}{\sqrt{1 + \alpha^2}} \quad \text{and} \quad b = |\mathbf{e}_2 \cdot \mathbf{v}| = \frac{|p - q\alpha|}{\sqrt{1 + \alpha^2}},$$

where  $p = m - k$  and  $q = n - \ell$ . Thus the area of the rectangle is equal to

$$\text{area} = ab = \frac{|p\alpha + q||p - q\alpha|}{1 + \alpha^2}. \quad (4.257)$$

Without loss of generality we can assume that  $p \geq 0$  and  $q \geq 0$ , and that  $p$  is the nearest integer to  $q\alpha$ . Then  $|p - q\alpha| = \|q\alpha\|$ . Next we need the following fact from the theory of continued fractions.

**Lemma 24.** *If  $1 \leq q < q_i$ , then*

$$\|q\alpha\| \geq \|q_{i-1}\alpha\| > \frac{1}{(a_i + 2)q_{i-1}}.$$

We postpone the proof of Lemma 24.

Now assume that

$$\frac{q_{i-1}}{4} \leq q < q_i. \quad (4.258)$$

Applying Lemma 24 and (4.258), we have

$$|p - q\alpha| = \|q\alpha\| \geq \|q_{i-1}\alpha\| > \frac{1}{(a_i + 2)q_{i-1}} \geq \frac{1}{4(a_i + 2)q}.$$

Substituting this in (4.257) and assuming (4.258), we have

$$\text{area} = ab = \frac{(p\alpha + q)|p - q\alpha|}{1 + \alpha^2} \geq \frac{q|p - q\alpha|}{1 + \alpha^2} \geq \frac{1}{4(a_i + 2)(1 + \alpha^2)}. \quad (4.259)$$

Let us elaborate on the meaning of (4.259). It is about a rectangle of slope  $1/\alpha$  which contains two lattice points  $P = (k, \ell)$  and  $Q = (m, n)$ ; in fact,  $P$  and  $Q$  are two vertices of the rectangle. We write the vector from  $P$  to  $Q$  as  $\mathbf{v} = (p, q)$  and, without loss of generality, we can assume that  $p \geq 0$  and  $q \geq 0$ , and that  $p$  is the nearest integer to  $q\alpha$ . If  $q$  is large, then  $\sqrt{1 + \alpha^2}q$  is very close to the diameter of this long and narrow rectangle. It means that  $q$  is basically a size parameter of the rectangle. Assume that the restriction (4.258) holds. Then the inequality (4.259) tells us that the area of this long and narrow rectangle is at least  $1/4(a_i + 2)(1 + \alpha^2)$ , that is, the area is not too small if  $a_i$  is not too large.

We can therefore rephrase (4.258) and (4.259) together in a nutshell as follows. A small digit  $a_i$  yields the rectangle property locally. This means that we have a good chance to adapt the Riesz product technique.

For the convenience of the reader, we interrupt the argument, and include a proof of Lemma 24 which is surprisingly tricky.

*Proof of Lemma 24.* Recall (4.244), that

$$p_k = a_k p_{k-1} + p_{k-2}, \quad q_k = a_k q_{k-1} + q_{k-2}.$$

These recurrences hold for any  $a_k$ , including arbitrary real values. Writing

$$\alpha = [a_1, \dots, a_{k-1}, \alpha_k],$$

with

$$\alpha_k = a_k + \frac{1}{a_{k+1} + \frac{1}{a_{k+2} + \dots}} = [a_k; a_{k+1}, a_{k+2}, \dots],$$

we obtain the useful formula

$$\alpha = \frac{\alpha_k p_{k-1} + p_{k-2}}{\alpha_k q_{k-1} + q_{k-2}},$$

and it follows that

$$q_{k-1}\alpha - p_{k-1} = \frac{q_{k-1}p_{k-2} - p_{k-1}q_{k-2}}{\alpha_k q_{k-1} + q_{k-2}}. \tag{4.260}$$

It is not difficult to show that

$$q_{k-1}p_{k-2} - p_{k-1}q_{k-2} = -(q_{k-2}p_{k-3} - p_{k-2}q_{k-3}). \tag{4.261}$$

Since  $p_0 = 0, q_0 = 1, p_1 = 1$  and  $q_1 = a_1$ , we have  $q_1 p_0 - p_1 q_0 = -1$ . It follows by induction, using (4.261), that

$$q_{k-1}p_{k-2} - p_{k-1}q_{k-2} = (-1)^{k-1}. \tag{4.262}$$

Combining this with (4.260), we have

$$q_{k-1}\alpha - p_{k-1} = \frac{(-1)^{k-1}}{\alpha_k q_{k-1} + q_{k-2}}, \tag{4.263}$$

which implies

$$\|q_{k-1}\alpha\| = |q_{k-1}\alpha - p_{k-1}| = \frac{1}{\alpha_k q_{k-1} + q_{k-2}} > \frac{1}{(a_k + 2)q_{k-1}}.$$

It remains to prove that, if  $p$  and  $q$  are integers with  $0 < q < q_k$ , then

$$|q\alpha - p| \geq |q_{k-1}\alpha - p_{k-1}|. \quad (4.264)$$

To prove this, we define integers  $u$  and  $v$  by the equations

$$p = up_{k-1} + vp_k, \quad q = uq_{k-1} + vq_k. \quad (4.265)$$

Note that (4.265) is solvable in integers  $u$  and  $v$ , since the determinant of the system is  $\pm 1$ , in view of (4.262). Since  $0 < q < q_k$ , we must have  $u \neq 0$ . Moreover, if  $v \neq 0$ , then  $u$  and  $v$  must have opposite signs. Since  $q_{k-1}\alpha - p_{k-1}$  and  $q_k\alpha - p_k$  also have opposite signs, in view of (4.263), we conclude that

$$\begin{aligned} |q\alpha - p| &= |u(q_{k-1}\alpha - p_{k-1}) + v(q_k\alpha - p_k)| \\ &= |u(q_{k-1}\alpha - p_{k-1})| + |v(q_k\alpha - p_k)| \\ &\geq |u(q_{k-1}\alpha - p_{k-1})| \geq |q_{k-1}\alpha - p_{k-1}|, \end{aligned}$$

proving (4.264). □

*Step 3.* We next employ the Riesz product technique. Let us return to Theorem 12, and the basically equivalent Proposition 13. A trivial novelty is that in this section, the slope is  $1/\alpha$ , whereas in Theorem 12 and Proposition 13, the slopes are respectively  $1/\sqrt{2}$  and 0. The Riesz product (4.99) is defined by using some appropriate modified Rademacher functions  $R_j(\mathbf{x}) \in \mathcal{R}(j)$  for  $j$  with  $1 \leq 2^j \leq N$ , i.e. for  $\log_2 N + O(1)$  values of  $j$ . In the hypothesis of Theorem 12 and Proposition 13, we have the unrestricted rectangle property; here we have a restricted rectangle property instead, meaning that the rectangle property holds only for  $O(\log N)$  values of the power-of-two parameter  $j$ , where  $0 \leq j \leq \log_2 N + O(1)$ . Indeed, by (4.250) and (4.251), we have

$$\log N_k = \log q_{k-1} + O(1) \leq \log \prod_{i=1}^k (a_i + 1) + O(1) = O(\log N),$$

and by (4.256), the continued fraction digit  $a_i$  of  $\alpha$  is small for at least  $k/2$  values of  $i$  in  $1 \leq i \leq k$ , if  $k$  is sufficiently large. For these small values of the continued fraction  $a_i$ , the rectangle property holds in the power-of-two range around  $q_{i-1}$ , i.e. when  $2^j \approx q_{i-1}$ ; see (4.258) and (4.259). This means that we can easily save the Riesz product technique developed earlier in Sects. 4.5–4.8. The minor price that we pay is a constant factor loss, due to the fact that  $\log_2 N$  is replaced by  $c_{28} \log N$ , where  $c_{28} = c_{28}(\gamma)$  is a small positive constant depending only on  $\gamma > 0$ . Thus we obtain the following result.

**Lemma 25.** *Let  $I = [a, b]$ , where  $0 \leq a < b < 1$ , be an arbitrary subinterval of the unit interval. Assume that (4.249) holds. Then there exists a constant  $\delta' = \delta'(\gamma) > 0$ , depending only on  $\gamma > 0$ , such that the following hold:*

- (i) *For all sufficiently large integers  $N$ , there is a subinterval  $I_1 = [a_1, b_1]$  of  $I$ , possibly depending on  $N$  and with  $a < a_1 < b_1 < b$ , such that for all  $\beta_1 \in I_1$ ,*

$$f(\alpha; \beta_1; \gamma; N) > 2\gamma \log N + \delta' \log N.$$

- (ii) *For all sufficiently large integers  $N$ , there is a subinterval  $I_2 = [a_2, b_2]$  of  $I$ , possibly depending on  $N$  and with  $a < a_2 < b_2 < b$ , such that for all  $\beta_2 \in I_2$ ,*

$$f(\alpha; \beta_2; \gamma; N) < 2\gamma \log N - \delta' \log N.$$

*Step 4.* The last step, the construction of a Cantor set, is routine. Combining the method of nested intervals with Lemma 25, we can easily build an infinite binary tree of nested intervals the same way as in the proof of Theorem 3. The divergence points  $\beta^*$  arise as the intersection of infinitely many decreasing intervals, which correspond to an infinite branch of the binary tree. Since a binary tree of countably infinite height has continuum many infinite branches, we obtain continuum many divergence points, proving Theorem 21 in Case 22.

*Case 26.* The inequality

$$\sum_{i=1}^k \lfloor \sqrt{\gamma a_i} \rfloor > 100 \cdot 2\gamma \log N_k \tag{4.266}$$

holds for infinitely many integers  $k \geq 1$ , where  $N_k$  is defined by (4.247).

The estimate (4.241) tells us that  $2\gamma \log N_k$  is the average value of  $f(\alpha; \beta; \gamma; N_k)$  as  $\beta$  runs through the unit interval. On the other hand, combining (4.248) and (4.266), we deduce that

$$f(\alpha; \beta = 0; \gamma; N_k) > 100 \cdot 2\gamma \log N_k$$

for infinitely many integers  $k \geq 1$ . In other words, for infinitely many values  $N = N_k$ , the homogeneous case  $\beta = 0$  gives at least 100 times more integer solutions than the average value  $2\gamma \log N_k$ . This represents an extreme bias; in fact, an extreme surplus. The proof of Theorem 3 is based on a somewhat similar extreme bias, a violation of the Naive Area Principle, in the sense that the Pell inequality  $-1 < x^2 - 2y^2 < 1$  has no integer solution except  $x = y = 0$ , while the corresponding hyperbolic region has infinite area. The only difference is that whereas in Theorem 3, we have an extreme shortage of solutions for the homogeneous case  $\beta = 0$ , we have here an extreme surplus. But this difference is irrelevant for the method of nested intervals, as it works in both cases. This means

that in Case 2, we can simply repeat the Cantor set construction in the proof of Theorem 3. This completes the proof of Theorem 21.

Theorem 21 is a qualitative result. In contrast, we complete this section with a quantitative result.

**Proposition 27.** *Let  $\alpha > 0$  and  $\gamma > 0$  be arbitrary real numbers. Then there is an effectively computable positive constant  $\delta' = \delta'(\gamma) > 0$ , depending only on  $\gamma > 0$ , such that for every sufficiently large integer  $N$ , there exist two real numbers  $\beta_1(N)$  and  $\beta_2(N)$  in the unit interval, with  $0 \leq \beta_1(N) < \beta_2(N) < 1$ , such that*

$$|f(\alpha; \beta_1(N); \gamma; N) - f(\alpha; \beta_2(N); \gamma; N)| > \delta' \log N.$$

We just outline the proof in a couple of sentences, since it is basically the same as that of Theorem 21, without the Cantor set construction. Indeed, let  $q_{\ell-1} \leq N < q_\ell$ . Since  $q_\ell = a_\ell q_{\ell-1} + q_{\ell-2} \leq (a_\ell + 1)q_{\ell-1}$ , we have

$$1 \leq \frac{N}{q_{\ell-1}} \leq a_\ell + 1.$$

Again we distinguish two cases.

Case 28. We have

$$\sum_{i=1}^{\ell-1} \lfloor \sqrt{\gamma a_i} \rfloor + \left\lfloor \sqrt{\frac{\gamma N}{q_{\ell-1}}} \right\rfloor \leq 100 \cdot 2\gamma \log N.$$

Then by repeating the argument of Case 1 in the proof of Theorem 21 above, we obtain Proposition 27; see Lemma 25.

Case 29. We have

$$\sum_{i=1}^{\ell-1} \lfloor \sqrt{\gamma a_i} \rfloor + \left\lfloor \sqrt{\frac{\gamma N}{q_{\ell-1}}} \right\rfloor > 100 \cdot 2\gamma \log N.$$

Then

$$f(\alpha; \beta = 0; \gamma; N) > 100 \cdot 2\gamma \log N,$$

and so we can choose  $\beta_1(N) = 0$ . Finally, for  $\beta_2(N)$ , we can choose any below average point; in other words, we can choose  $\beta_2(N)$  to be any  $\beta$  that satisfies the inequality  $f(\alpha; \beta; \gamma; N) \leq (2 + o(1))\gamma \log N$ ; see (4.241).



## 4.10 General Point Sets: Theorem 30

What will happen if we drop the rectangle property in Theorem 12 or Proposition 13? Can we still exhibit extra large deviations for hyperbolic needles? This is the subject of this last section.

Suppose that  $\mathcal{P}$  is a finite point set of density  $\delta > 0$  in a large square  $[0, M]^2$ , so that  $|\mathcal{P}| = \delta M^2$ . We shall make a very mild technical assumption, that  $\mathcal{P}$  is not clustered. More precisely, we introduce a new concept called the *separation constant* and denoted by  $\sigma = \sigma(\mathcal{P})$ , and say that  $\mathcal{P}$  is  $\sigma$ -separated if the usual Euclidean distance between any two points of  $\mathcal{P}$  is at least  $\sigma$ . For example, the set of integer lattice points in the plane is clearly 1-separated, so that  $\sigma(\mathbf{Z}^2) = 1$ .

Our basic idea is the following. We show that if  $\mathcal{P}$  is  $\sigma$ -separated with some not too small constant  $\sigma > 0$ , then the rectangle property holds, at least in a weak statistical sense, for the majority of the directions which we shall call the good directions. For example, in Theorem 12, the slope  $1/\sqrt{2}$  is a concrete good direction. This is how we will be able to save the Riesz product argument in the proof of Theorem 12 or Proposition 13, and still prove extra large deviations, proportional to the area, for hyperbolic needles, at least for the majority of the directions.

In the rest of the section, we work out the details of the vague intuition, and this will give us Theorem 30. The obvious handicap of this majority approach is that for an arbitrary point set  $\mathcal{P}$  which is not clustered, we cannot predict whether a given concrete direction is good or not.

Another, and purely technical, shortcoming is that in Theorem 30, we cannot get rid of the assumption that  $\mathcal{P}$  is not clustered. This technical difficulty is rather counterintuitive, since at least at first sight, clusters actually seem to help us create extra large deviations. However, some technical difficulties prevent us from adapting the Riesz product technique for clustered point sets  $\mathcal{P}$ . It remains an interesting open problem to decide whether or not the separation constant  $\sigma = \sigma(\mathcal{P})$  in Theorem 30 plays any role.

In Theorem 30, we change<sup>15</sup> the underlying set, and switch from the large square  $[0, M]^2$  to the large disk

$$\text{disk}(\mathbf{0}; M) = \{\mathbf{x} \in \mathbf{R}^2 : |\mathbf{x}| \leq M\}$$

of radius  $M$  and centered at the origin.

Let  $\mathcal{P}$  be a finite point set of density  $\delta > 0$  in the large disk  $\text{disk}(\mathbf{0}; M)$ , so that  $|\mathcal{P}| = \delta\pi M^2$ ; here we assume that the radius  $M$  is large. We also assume that  $\mathcal{P}$  is not clustered. More precisely, we assume that  $\mathcal{P}$  is  $\sigma$ -separated for some positive constant  $\sigma = \sigma(\mathcal{P}) > 0$ . The goal is to count the number of elements of  $\mathcal{P}$  in rotated and translated copies of our usual hyperbolic needle  $H_\gamma(N)$ .

---

<sup>15</sup>The reason behind this change is rotation-invariance. Theorems 3 and 12 are about translated copies, whereas Theorem 30 is about rotated and translated copies of the hyperbolic needle.

Let  $10^{-2} > \eta > 0$  be a small positive real numbers, to be specified later. Let  $j$  be an arbitrary integer in the interval  $0 \leq j \leq n$ , where  $2^n \approx N$ , that is,  $n = \log_2 N + O(1)$  in binary logarithm. We decompose the large disk  $\text{disk}(\mathbf{0}; M)$  into disjoint translated copies of the small rectangle

$$[0, 2^j \eta] \times [0, 2^{-j} \eta]; \tag{4.267}$$

in other words, we form a rectangle lattice starting from the origin. We shall focus on the copies of (4.267) which are inside the large disk  $\text{disk}(\mathbf{0}; M)$ , and ignore the copies of (4.267) that intersect the boundary circle or are outside the disk. Note that there are  $O(2^j \eta M)$  copies of (4.267) that intersect the boundary circle of the large disk. If  $2^j \eta = o(M)$ , then there are  $(1 + o(1))\pi M^2 \eta^{-2}$  copies of (4.267) that are inside the large disk  $\text{disk}(\mathbf{0}; M)$ . We call these translated copies of the small rectangle (4.267)  $j$ -cells. More precisely, we call them  $j$ -cells of angle 0.

In general, let  $\theta$  be an arbitrary angle, with  $0 \leq \theta < \pi$ . Let  $\text{Rot}_\theta$  denote the rotation of the plane by the angle  $\theta$ , assuming that the fixed point of the rotation  $\text{Rot}_\theta$  is the origin. We decompose the large disk  $\text{disk}(\mathbf{0}; M)$  into disjoint translates of the rotated copy

$$\text{Rot}_\theta([0, 2^j \eta] \times [0, 2^{-j} \eta]) \tag{4.268}$$

of the small rectangle (4.267). We shall focus on the translated copies of (4.268) which are inside the large disk  $\text{disk}(\mathbf{0}; M)$ . Again, if  $2^j \eta = o(M)$ , then there are  $(1 + o(1))\pi M^2 \eta^{-2}$  translated copies of (4.268) that are inside the large disk  $\text{disk}(\mathbf{0}; M)$ . We call these translated copies of the small rectangle (4.268)  $j$ -cells of angle  $\theta$ .

We want to prove, in a quantitative form, that if  $\mathcal{P}$  is not clustered, then for a typical angle  $\theta \in [0, \pi)$ , the overwhelming majority of the  $j$ -cells of angle  $\theta$  that contain at least one point of  $\mathcal{P}$  actually contain exactly one point of  $\mathcal{P}$ . A quantitative result like this, a statistical version of the rectangle property, will serve as a substitute for the rectangle property, and it will suffice to save the Riesz product technique developed in Sects. 4.5–4.8.

*Statistical Version of the Rectangle Property: An Average Argument.* Suppose that  $P_{i_1}, P_{i_2} \in \mathcal{P}$ , where  $i_1 \neq i_2$ , are two arbitrary points. We define the *angle-set* by

$$\text{angle}(P_{i_1}, P_{i_2}; j) = \{\theta \in [0, \pi) : \text{there is a } j\text{-cell of angle } \theta \text{ containing } P_{i_1} \text{ and } P_{i_2}\}.$$

The angle-set  $\text{angle}(P_{i_1}, P_{i_2}; j)$  is clearly measurable. Let  $|\text{angle}(P_{i_1}, P_{i_2}; j)|$  denote the usual one-dimensional Lebesgue measure, i.e. length.

The basic idea is to estimate the double sum

$$\sum_{\substack{P_{i_1}, P_{i_2} \in \mathcal{P} \\ i_1 \neq i_2}} |\text{angle}(P_{i_1}, P_{i_2}; j)|.$$

Simple geometric consideration shows that

$$|\text{angle}(P_{i_1}, P_{i_2}; j)| < 2 \cdot \frac{2^{-j} \eta}{|P_{i_1} P_{i_2}|},$$

where  $2^{-j} \eta$  is the length of the short side of a  $j$ -cell and  $|P_{i_1} P_{i_2}|$  denotes the usual Euclidean distance between  $P_{i_1}$  and  $P_{i_2}$ , and so

$$\sum_{\substack{P_{i_1}, P_{i_2} \in \mathcal{P} \\ i_1 \neq i_2}} |\text{angle}(P_{i_1}, P_{i_2}; j)| < 2^{-j} \eta \sum_{P_{i_1} \in \mathcal{P}} \left( \sum_{\substack{P_{i_2} \in \mathcal{P} \\ i_1 \neq i_2}} \frac{1}{|P_{i_1} P_{i_2}|} \right). \quad (4.269)$$

Since  $\mathcal{P}$  is  $\sigma$ -separated, it is easy to give an upper bound to the inner sum in (4.269). Using a standard power-of-two decomposition, we have

$$\begin{aligned} \sum_{\substack{P_{i_2} \in \mathcal{P} \\ i_1 \neq i_2}} \frac{1}{|P_{i_1} P_{i_2}|} &\leq \sum_{1 \leq \ell \leq L} \sum_{\substack{P_{i_2} \in \mathcal{P} \\ i_1 \neq i_2 \\ 2^{\ell-1} \sigma < |P_{i_1} P_{i_2}| \leq 2^\ell \sigma}} \frac{1}{|P_{i_1} P_{i_2}|} \\ &\leq \sum_{1 \leq \ell \leq L} \frac{1}{2^{\ell-1} \sigma} \cdot \pi(2^{\ell+1})^2 = \sum_{1 \leq \ell \leq L} \frac{8\pi}{\sigma} \cdot 2^\ell < \frac{16\pi}{\sigma} \cdot 2^L, \end{aligned} \quad (4.270)$$

where  $L$  denotes the largest integer such that  $2^L \sigma < 2^{j+1} \eta$ , and where the estimate  $\pi(2^{\ell+1})^2$  arises from the fact that a square of side  $\sigma/2$  cannot contain two points from  $\mathcal{P}$ , since  $\mathcal{P}$  is  $\sigma$ -separated. From (4.270), and using the fact that  $2^L \sigma < 2^{j+1} \eta$ , we conclude that

$$\sum_{\substack{P_{i_2} \in \mathcal{P} \\ i_1 \neq i_2}} \frac{1}{|P_{i_1} P_{i_2}|} < \frac{16\pi}{\sigma} \cdot 2^L < \frac{16\pi}{\sigma} \cdot \frac{2^{j+1} \eta}{\sigma} = \frac{2^5 \pi \eta 2^j}{\sigma^2}. \quad (4.271)$$

Combining (4.269) and (4.271), and using the fact that  $|\mathcal{P}| = \delta \pi M^2$ , we then obtain

$$\sum_{\substack{P_{i_1}, P_{i_2} \in \mathcal{P} \\ i_1 \neq i_2}} |\text{angle}(P_{i_1}, P_{i_2}; j)| < 2^{-j} \eta |\mathcal{P}| \frac{2^5 \pi \eta 2^j}{\sigma^2} = \frac{2^5 \pi^2 \eta^2 \delta M^2}{\sigma^2}. \quad (4.272)$$

---

<sup>16</sup>Note that  $2^j \eta$  is the length of the long side of a  $j$ -cell.

Recall that the disk  $\text{disk}(\mathbf{0}; M)$  of radius  $M$  contains  $(1 + o(1))\pi M^2 \eta^{-2}$   $j$ -cells of a given angle  $\theta$ , and that  $\theta$  runs through the interval  $0 \leq \theta < \pi$ . It is natural, therefore, to normalize the sum (4.272) and consider the average

$$\frac{1}{\pi^2 M^2 \eta^{-2}} \sum_{\substack{P_{i_1}, P_{i_2} \in \mathcal{P} \\ i_1 \neq i_2}} |\text{angle}(P_{i_1}, P_{i_2}; j)| < \eta^4 \cdot \frac{2^5 \delta}{\sigma^2}. \tag{4.273}$$

*Consequences of Inequality (4.273).* Let us return to Sect. 4.8. Recall that the last step in the proof of Proposition 13, and indirectly the proof of Theorem 12, is to choose the parameters  $\eta_1$  and  $\eta_2$  as sufficiently small positive constants independent of  $M$  and  $N$ ; see (4.226). In fact, in view of (4.100),  $\eta_1$  and  $\eta_2$  are almost equal.

In similar fashion, we assume here that the parameter  $\gamma$  of the hyperbolic needle, the density  $\delta$  of  $\mathcal{P}$  and the separation constant  $\sigma$  of  $\mathcal{P}$  are fixed positive constants, and consider  $\eta$ , which of course plays the role of  $\eta_1$  and  $\eta_2$ , as a parameter that we shall eventually choose as a sufficiently small positive constant independent of  $M$  and  $N$ .

Since the area of a  $j$ -cell is  $\eta^2$ , we can say roughly that the probability that a  $j$ -cell of any angle contains a point of  $\mathcal{P}$  is

$$\text{density} \times \text{area} = \delta \eta^2. \tag{4.274}$$

On the other hand, in view of (4.273), the probability that a  $j$ -cell of any angle contains exactly two points of  $\mathcal{P}$  does not exceed  $c_{29} \eta^4$ , which is negligible compared to  $\delta \eta^2$  in (4.274) if  $\eta$  is small enough.

In general, the probability that a  $j$ -cell of any angle contains exactly  $p$  points of  $\mathcal{P}$ , where  $2^\ell < p \leq 2^{\ell+1}$  with  $\ell = 1, 2, 3, \dots$ , does not exceed  $c_{30} \eta^4 4^{-\ell}$ , where the constant factor  $c_{30}$  is independent of  $\ell$ . Indeed,  $p$  points from  $\mathcal{P}$  means that we can choose  $\binom{p}{2}$  pairs  $P_{i_1}, P_{i_2}$ , implying that those rich  $j$ -cells show up with multiplicity

$$\binom{p}{2} > 2^\ell 2^{\ell-1} = \frac{1}{2} 4^\ell$$

in (4.273), explaining the factor  $4^{-\ell}$  in  $c_{30} \eta^4 4^{-\ell}$ . The point here is that even the sum of the products

$$\sum_{\ell \geq 1} 2^{\ell+1} \eta^4 4^{-\ell}$$

is negligible compared to the  $\delta \eta^2$  in (4.274) if  $\eta$  is small enough.

Summarizing, we can say that (4.273) implies the following general picture about the distribution of the elements of  $\mathcal{P}$  in the  $j$ -cells of any angle. Let  $\theta \in [0, \pi)$  be a typical angle, and consider the  $j$ -cells of angle  $\theta$ . The overwhelming majority of

the points  $P \in \mathcal{P}$  turn out to be singles, meaning that if the point  $P$  is contained in some  $j$ -cell  $\mathcal{C}$  of angle  $\theta$ , then  $\mathcal{C}$  does not contain any other point of  $\mathcal{P}$ . Here the vague term *overwhelming majority* in fact has the quantitative meaning of  $1 - O(\eta^2)$  part of  $\mathcal{P}$ . Note that  $1 - O(\eta^2)$  is almost 1 if  $\eta$  is small.

Furthermore, rich  $j$ -cells turn out to be very rare in the following sense. Let  $\ell \geq 0$  be a fixed integer. The proportion of the  $j$ -cells  $\mathcal{C}$  of angle  $\theta$  containing  $p$  points of  $\mathcal{P}$ , where  $2^\ell < p \leq 2^{\ell+1}$ , compared to those  $j$ -cells which contain at least one point of  $\mathcal{P}$ , does not exceed  $c_{31}\eta^2 4^{-\ell}$ , where the constant factor  $c_{31}$  is independent of  $\ell$ . Since  $2^\ell$  is negligible compared to  $4^\ell$  if  $\ell$  is large, the term *very rare* is well justified.

We can say, therefore, that a weaker statistical version of the rectangle property holds for the majority of the angles  $\theta \in [0, \pi)$ , assuming that  $\eta > 0$  is a sufficiently small constant depending only on the parameter  $\gamma$  of the hyperbolic needle, the density  $\delta$  of  $\mathcal{P}$  and the separation constant  $\sigma$  of  $\mathcal{P}$ .

A simple analysis of the Riesz product argument in Sects. 4.5–4.8 shows that this weaker statistical version of the rectangle property is a good substitute for the strict rectangle property, and thus we can prove the following result.

**Theorem 30.** *Let  $\mathcal{P}$  be a finite set of points in the disk  $\text{disk}(\mathbf{0}; M)$  with density  $\delta$ , so that the number of elements of  $\mathcal{P}$  is  $|\mathcal{P}| = \delta\pi M^2$ . Assume that  $\mathcal{P}$  is  $\sigma$ -separated for some  $\sigma > 0$ . Assume further that both  $N$  and  $M/N$  are sufficiently large, depending only on  $\gamma$ ,  $\delta$  and  $\sigma$ . Then there exist a positive constant  $\delta' = \delta'(\sigma, \gamma, \delta) > 0$ , independent of  $N$  and  $M$ , and a measurable subset  $\mathcal{A} \subset [0, 2\pi)$ , of Lebesgue measure greater than  $\frac{99}{100} \cdot 2\pi$ , such that for every angle  $\theta \in \mathcal{A}$ , there exist translated copies  $\mathbf{x}_1 + \text{Rot}_\theta H_\gamma(N) \subset \text{disk}(\mathbf{0}; M)$  and  $\mathbf{x}_2 + \text{Rot}_\theta H_\gamma(N) \subset \text{disk}(\mathbf{0}; M)$  of the rotated hyperbolic needle  $\text{Rot}_\theta H_\gamma(N)$  such that*

$$|\mathcal{P} \cap (\mathbf{x}_1 + \text{Rot}_\theta H_\gamma(N))| \geq 2\delta\gamma \log N + \delta' \log N$$

and

$$|\mathcal{P} \cap (\mathbf{x}_2 + \text{Rot}_\theta H_\gamma(N))| \leq 2\delta\gamma \log N - \delta' \log N.$$

As indicated at the beginning of this section, it is reasonable to guess that clusters just help to create extra large fluctuations. This intuition motivates the following

**Open Problem.** *Can one prove a version of Theorem 30 which makes no reference to the separation constant  $\sigma = \sigma(\mathcal{P})$ ? In other words, can we simply drop  $\sigma = \sigma(\mathcal{P})$  from the hypotheses of Theorem 30?*

The author guesses that the answer is affirmative but, unfortunately, cannot prove it.

Finally, we briefly mention a closely related problem, where we cannot drop the separation constant  $\sigma = \sigma(\mathcal{P})$  from the hypotheses. Note that Theorems 3–21 all concern the extra large fluctuations of the measure-theoretic discrepancy, meaning the difference between the number of points of  $\mathcal{P}$  and its expectation of density

times area. What we study last here is the large fluctuations of the  $\pm 1$ -discrepancy, or 2-coloring discrepancy.

This means that we have an arbitrary 2-coloring  $\varphi : \mathcal{P} \rightarrow \{\pm 1\}$  of the given point set  $\mathcal{P}$ , with  $+1$  representing *red* and  $-1$  representing *blue*, say. Extra large fluctuations of the  $\pm 1$ -discrepancy means that there is a translated, or rotated and translated, copies  $H'$  and  $H''$  of the hyperbolic needle  $H_\gamma(N)$  such that

$$\sum_{P \in \mathcal{P} \cap H'} \varphi(P) > c_{32} \cdot \text{area}(H') = c_{33} \log N > 0$$

with some positive constants  $c_{32}$  and  $c_{33}$  and

$$\sum_{P \in \mathcal{P} \cap H''} \varphi(P) < -c_{34} \cdot \text{area}(H'') = -c_{35} \log N < 0$$

with some positive constants  $c_{34}$  and  $c_{35}$ .

The Riesz product technique can be easily adapted to prove extra large fluctuations of the  $\pm 1$ -discrepancy. For example, we have the following  $\pm 1$ -discrepancy analog of Proposition 13.

**Proposition 31 (2-Coloring Discrepancy for Translated Copies).** *Let  $\mathcal{P}$  be a finite set of points in the square  $[0, M]^2$  with density  $\delta$ , so that the number of elements of  $\mathcal{P}$  is  $|\mathcal{P}| = \delta M^2$ . Let  $\varphi : \mathcal{P} \rightarrow \{\pm 1\}$  be an arbitrary 2-coloring of  $\mathcal{P}$ . Assume that  $\mathcal{P}$  satisfies the following rectangle property, that there is a positive constant  $c_1 = c_1(\mathcal{P}) > 0$  such that every axes-parallel rectangle of area  $c_1$  contains at most one element of the set  $\mathcal{P}$ . As in Proposition 13, let  $\delta' = \delta'(c_1, \gamma, \delta)$  be defined by (4.83) and (4.84), and assume that both  $N$  and  $M/N$  are sufficiently large and satisfy (4.85). Then for the hyperbolic needle  $H_\gamma(N)$  given by (4.86), there exist translated copies  $\mathbf{x}_1 + H_\gamma(N) \subset [0, M]^2$  and  $\mathbf{x}_2 + H_\gamma(N) \subset [0, M]^2$  such that*

$$\sum_{P \in \mathcal{P} \cap (\mathbf{x}_1 + H_\gamma(N))} \varphi(P) \geq \delta' \log N$$

and

$$\sum_{P \in \mathcal{P} \cap (\mathbf{x}_2 + H_\gamma(N))} \varphi(P) \leq -\delta' \log N.$$

Similarly, one can easily prove the following analog of Theorem 30.

**Proposition 32 (2-Coloring Discrepancy for Rotated and Translated Copies).** *Let  $\mathcal{P}$  be a finite set of points in the disk  $\text{disk}(\mathbf{0}; M)$  with density  $\delta$ , so that the number of elements of  $\mathcal{P}$  is  $|\mathcal{P}| = \delta \pi M^2$ . Let  $\varphi : \mathcal{P} \rightarrow \{\pm 1\}$  be an arbitrary 2-coloring of  $\mathcal{P}$ . Assume that  $\mathcal{P}$  is  $\sigma$ -separated with some  $\sigma > 0$ . Assume further that both  $N$  and  $M/N$  are sufficiently large, depending only on  $\gamma$ ,  $\delta$  and  $\sigma$ . Then*

there exist a positive constant  $\delta' = \delta'(\sigma, \gamma, \delta) > 0$ , independent of  $N$  and  $M$ , and a measurable subset  $\mathcal{A} \subset [0, 2\pi)$ , of Lebesgue measure greater than  $\frac{99}{100} \cdot 2\pi$ , such that for every angle  $\theta \in \mathcal{A}$ , there exist translated copies  $\mathbf{x}_1 + \text{Rot}_\theta H_\gamma(N) \subset \text{disk}(\mathbf{0}; M)$  and  $\mathbf{x}_2 + \text{Rot}_\theta H_\gamma(N) \subset \text{disk}(\mathbf{0}; M)$  of the rotated hyperbolic needle  $\text{Rot}_\theta H_\gamma(N)$  such that

$$\sum_{P \in \mathcal{P} \cap (\mathbf{x}_1 + \text{Rot}_\theta H_\gamma(N))} \varphi(P) \geq \delta' \log N$$

and

$$\sum_{P \in \mathcal{P} \cap (\mathbf{x}_2 + \text{Rot}_\theta H_\gamma(N))} \varphi(P) \leq -\delta' \log N.$$

We want to point out that in Proposition 32 on the  $\pm 1$ -discrepancy of hyperbolic needles, we definitely need some extra condition implying that  $\mathcal{P}$  is not too clustered. Indeed, it is easy to construct an extremely clustered point set  $\mathcal{P}$  for which the  $\pm 1$ -discrepancy of the hyperbolic needles is negligible. For example, we can start with a typical point set in general position, and split up every point into a pair of points being extremely close to each other. The two points in these extremely close pairs are joined with a straight line segment each, and we refer to these line segments as the very short line segments. Consider the particular 2-coloring of the point set where the two points in the extremely close pairs all have different colors, with one  $+1$  and the other  $-1$ . We can easily guarantee that this particular 2-coloring has negligible  $\pm 1$ -discrepancy for the family of all hyperbolic needles congruent to  $H_\gamma(N)$ . If the original point set is in general position and the point pairs are close enough, then the arcs of any congruent copy of  $H_\gamma(N)$  intersect at most two very short line segments. Since the boundary of  $H_\gamma(N)$  consists of 4 arcs, the  $\pm 1$ -discrepancy is at most  $4 \cdot 2 = 8$ , which is indeed negligible.

## References

1. J. Beck, Randomness of  $n\sqrt{2} \bmod 1$  and a Ramsey property of the hyperbola, in *Sets, Graphs and Numbers*, ed. by G. Halász, L. Lovász, D. Miklós, T. Szőnyi. Colloquia Math. Soc. János Bolyai, vol. 60 (North-Holland Publishing, Amsterdam, 1992), pp. 23–66
2. J. Beck, *Randomness in diophantine approximation* (Springer, to appear)
3. J. Beck, Diophantine approximation and quadratic fields, in *Number Theory*, ed. by K. Györy, A. Pethő, V.T. Sós. (Walter de Gruyter, Berlin, 1998), pp. 53–93
4. J. Beck, From probabilistic diophantine approximation to quadratic fields, in *Random and Quasi-Random Point Sets*, ed. by Hellekalek, P., Larcher G. Lecture Notes in Statistics, vol. 138 (Springer, New York, NY, 1998), pp. 1–48
5. J. Beck, Randomness in lattice point problems. *Discrete Math.* **229**(1–3), 29–55 (2001). doi:[10.1016/S0012-365X\(00\)00200-4](https://doi.org/10.1016/S0012-365X(00)00200-4)

6. J. Beck, Lattice point problems: crossroads of number theory, probability theory and Fourier analysis, in *Fourier Analysis and Convexity*, ed. by L. Brandolini, L. Colzani, A. Iosevich, G. Travaglini (Birkhäuser, Boston, MA, 2004), pp. 1–35
7. J. Beck, *Inevitable randomness in discrete mathematics*. University Lecture Series, vol. 49 (American Mathematical Society (AMS), Providence, RI, 2009)
8. J. Beck, Lattice point counting and the probabilistic method. *J. Combinator.* **1**(2), 171–232 (2010)
9. J. Beck, Randomness of the square root of 2 and the giant leap. I. *Period. Math. Hung.* **60**(2), 137–242 (2010). doi:[10.1007/s10998-010-2137-9](https://doi.org/10.1007/s10998-010-2137-9)
10. J. Beck, Randomness of the square root of 2 and the giant leap. II. *Period. Math. Hung.* **62**(2), 127–246 (2011). doi:[10.1007/s10998-011-6127-3](https://doi.org/10.1007/s10998-011-6127-3)
11. J. Beck, W.W.L. Chen, *Irregularities of distribution*. Cambridge Tracts in Mathematics vol. 89 (Cambridge University Press, Cambridge, 1987)
12. J.W.S. Cassels, An extension of the law of the iterated logarithm. *Math. Proc. Cambridge Philos. Soc.* **47**, 55–64 (1951)
13. B. Chazelle, *The discrepancy method. Randomness and complexity* (Cambridge University Press, Cambridge, 2000)
14. P. Erdős, On the law of the iterated logarithm. *Ann. Math. (2)* **43**, 419–436 (1942). doi:[10.2307/1968801](https://doi.org/10.2307/1968801)
15. W. Feller, The general form of the so-called law of the iterated logarithm. *Trans. Am. Math. Soc.* **54**, 373–402 (1943). doi:[10.2307/1990253](https://doi.org/10.2307/1990253)
16. W. Feller, *An introduction to probability theory and its applications. I.*, 3rd edn. (Wiley, New York, 1968)
17. W. Feller, *An introduction to probability theory and its applications. II.*, 2nd edn. (Wiley, New York, 1971)
18. R.L. Graham, B.L. Rothschild, J.H. Spencer, *Ramsey theory* (Wiley, New York, 1980)
19. G. Halász, On Roth’s method in the theory of irregularities of point distributions, in *Recent Progress in Analytic Number Theory*, vol. 2, ed. by H. Halberstam, C. Hooley, vol. 2 (Academic Press, London, 1981), pp. 79–94
20. M. Kac, Probability methods in some problems of analysis and number theory. *Bull. Am. Math. Soc.* **55**, 641–665 (1949). doi:[10.1090/S0002-9904-1949-09242-X](https://doi.org/10.1090/S0002-9904-1949-09242-X)
21. A. Khintchine, Über einen Satz der Wahrscheinlichkeitsrechnung. *Fund. math.* **6**, 9–20 (1924)
22. A. Kolmogorov, Über das Gesetz des iterierten Logarithmus. *Math. Ann.* **101**, 126–135 (1929). doi:[10.1007/BF01454828](https://doi.org/10.1007/BF01454828)
23. S. Lang, *Introduction to diophantine approximations* (Addison-Wesley, Reading, 1966)
24. J. Matoušek, *Geometric discrepancy. An illustrated guide. Revised paperback reprint of the 1999 original*. Algorithms and Combinatorics, vol. 18 (Springer, Berlin, 2010). doi:[10.1007/978-3-642-03942-3](https://doi.org/10.1007/978-3-642-03942-3)
25. A. Ostrowski, Bemerkungen zur Theorie der diophantischen Approximationen I. *Abh. Math. Sem. Univ. Hamburg* **1**, 77–98 (1921). doi:[10.1007/BF02940581](https://doi.org/10.1007/BF02940581)
26. K.F. Roth, On irregularities of distribution. *Mathematika* **1**, 73–79 (1954). doi:[10.1112/S0025579300000541](https://doi.org/10.1112/S0025579300000541)
27. W.M. Schmidt, Irregularities of distribution. VII. *Acta Arith.* **21**, 45–50 (1972)
28. J.G. van der Corput, Verteilungsfunktionen. I. *Proc. Kon. Ned. Akad. v. Wetensch. Amsterdam* **38**, 813–821 (1935)
29. J.G. van der Corput, Verteilungsfunktionen. II. *Proc. Kon. Ned. Akad. v. Wetensch. Amsterdam* **38**, 1058–1066 (1935)
30. H. Weyl, Über die Gleichverteilung von Zahlen mod. Eins. *Math. Ann.* **77**, 313–352 (1916). doi:[10.1007/BF01475864](https://doi.org/10.1007/BF01475864)
31. D.B. Zagier, *Zetafunktionen und quadratische Körper. Eine Einführung in die höhere Zahlentheorie* (Springer, Berlin, 1981)