# A localization theorem in the theory of diophantine approximation and an application to Pell's equation

by

A. Rockett and P. Szüsz (Stony Brook, N.Y.)

**Introduction.** Some thirty years ago the second of us showed (Szüsz [3]) that for any irrational number $\alpha$ the inequality

$$\|\alpha x\| < x^{\delta-1} \quad (0 < \delta < 1)$$

has a solution in any interval $(n, n^{1/\delta})$ where $\|z\|$ denotes the distance from $z$ to the nearest integer. Since, as is well known, $\|\alpha x\| < x^{-1}$ is solvable with natural numbers $x$, it would be interesting to give a "localization" for the numbers $x$ satisfying $\|\alpha x\| < x^{-1}$ but it is easy to see that one can give a counterexample to any such statement. Further, one could ask for a characterization of the natural numbers $x$ for which

$$(1) \qquad \qquad \|\alpha x\| < K/x$$

holds where $K$ is a positive constant.

Let the regular continued fraction representation of $\alpha$ be $[a_0; a_1, a_2, \ldots]$ and let the denominators of the convergents be $B_0 = 1, B_1, B_2, \ldots$ In Section 1 we prove

THEOREM 1. *Let* $K \geqslant 1/\sqrt{5}$ *and let* $\alpha$ *be an irrational number. Then all positive integers* $x$ *satisfying* (1) *have the form*

$$(2) \qquad \qquad c_{n+1}B_n + c_{n+2}B_{n+1} + \ldots + c_{n+m}B_{n+m-1}$$

*with* $m < C\log(2K+1)+3$, *where* $C$ *is an absolute constant and the coefficients* $c_{k+1}$ *satisfy* $0 \leqslant c_i < a_i$, $0 \leqslant c_{k+1} \leqslant a_{k+1}$ *for* $k > 0$ *and if* $c_{k+1} = a_{k+1}$ *then* $c_k = 0$.

This result extends the classical result of Legendre that $m = 1$ for $K \leqslant 1/2$ and that $m \leqslant 2$ for $K = 1$ (see Perron [2], Sections 13 and 16). In Section 2 we apply our result to Pell's equation $x^2 - dy^2 = N$. While this equation has been treated for $|N| < \sqrt{d}$ (see Perron [2] for references), we can drop this restriction and thus generalize the classical results.

**1. Proof of the main theorem.** We use the notations of Perron [2] for regular continued fractions:

$$\alpha = [a_0; a_1, a_2, \ldots],$$

$$[a_0; a_1, \ldots, a_k] = A_k/B_k \quad \text{where} \quad (A_k, B_k) = 1,$$

$$\zeta_k = [a_k; a_{k+1}, \ldots].$$

We set

$$D_k = B_k \alpha - A_k = (-1)^k/(B_k \zeta_{k+1} + B_{k-1}).$$

Since $A_{k+1} = a_{k+1} A_k + A_{k-1}$ and $B_{k+1} = a_{k+1} B_k + B_{k-1}$, it follows that $D_{k+1} = a_{k+1} D_k + D_{k-1}$. Since $D_{k+1} = -D_k/\zeta_{k+2}$, the $D_k$'s alternate in sign and their absolute values decrease monotonically to zero.

LEMMA 1.1 (Ostrowski [1]). *Every positive integer $x$ has a unique representation as*

(3)
$$x = \sum_{k=0}^{N} c_{k+1} B_k$$

*where $0 \le c_1 < a_1, 0 \le c_{k+1} \le a_{k+1}$ for $k > 0$ and if $c_{k+1} = a_{k+1}$ then $c_k = 0$.*

The proof can be done by induction on $x$.

Let $c_{n+1}$ be the first nonzero coefficient in the representation (3) of $x$ so that $c_{n+1} > 0$ and $c_{k+1} = 0$ for $0 \le k < n \le N$.

LEMMA 1.2. *We have*

$$|(c_{n+1}-1) D_n - D_{n+1}| < \left| \sum_{k=0}^{N} c_{k+1} D_k \right| < |c_{n+1} D_n - D_{n+1}|.$$

Proof. Since the $D_k$'s alternate in sign,

$$\left| \sum_{k=0}^{N} c_{k+1} D_k \right| > |c_{n+1} D_n + (a_{n+2}-1) D_{n+1} + a_{n+4} D_{n+3} + \ldots|$$

and the lower estimate follows since $a_{k+1} D_k = D_{k+1} - D_{k-1}$. The upper estimate is obtained similarly by considering $|c_{n+1} D_n + a_{n+3} D_{n+2} + a_{n+5} D_{n+4} + \ldots|$.

LEMMA 1.3. *For any integer $x > 1$ either*

(4)
$$\|\alpha x\| = \left| \sum_{k=0}^{N} c_{k+1} D_k \right|$$

*or $\|\alpha x\| > D_2$.*

Proof. A simple calculation shows that if $c_1 = c_2 = 0$ then the right-hand side of (4) is $< 1/2$. The exceptional cases occur when $a_1 = 1$ and $c_2 > 0$ and when $c_1 > 0$.

We now prove Theorem 1. From Lemmas 1.2 and 1.3 we see that $\|\alpha x\|$ is minimized for integers of the form

$$x = B_n + (a_{n+2}-1) B_{n+1} + a_{n+4} B_{n+3} + \ldots + a_{n+2m} B_{n+2m-1} = B_{n+2m} - B_{n+1}$$

and for such numbers we have that

$$\|\alpha x\| x > |-D_{n+1}| x > (B_{n+2m} - B_{n+1})/(B_{n+2} + B_{n+1}) > (G^{2m-3} - 1)/2$$

where $G = (1 + \sqrt{5})/2$ since $B_{k+2m}/B_{k+2} > G^{2m-3}$ for any $\alpha$. Thus if $\|\alpha x\| x < K$ we must have $(G^{2m-3} - 1)/2 < K$ and our result follows.

We note that an upper estimate for $\|\alpha x\| x$ would require an upper bound for the ratios $B_{k+1}/B_k$ and these ratios are unbounded for almost all $\alpha$.

**2. An application to Pell's equation.** We now consider the positive integer solutions of $x^2 - dy^2 = N$ where the integer $d > 1$ is not a perfect square. For $N > 0$, we have that $\|y\sqrt{d}\|^2 + 2\sqrt{d}\|y\sqrt{d}\| y = N$ and so $\|y\sqrt{d}\| y < N/2\sqrt{d}$. If $0 < N < \sqrt{d}$ then $\|y\sqrt{d}\| y < 1/2$ and we obtain the classical result that the only solutions are those given by the convergents of $\sqrt{d}$ (the case $N < 0$ follows in a similar manner since the convergents of $\sqrt{d}$ and $1/\sqrt{d}$ coincide with one trivial exception). For $|N| > \sqrt{d}$, it follows that the solutions will be as described by Theorem 1 (together with the corresponding sums of $A_k$'s).

With the notations of Perron [2] for the regular continued fraction of $\sqrt{d}$, we have

$$\zeta_0 = (\sqrt{d} + 0)/1, \quad \ldots, \quad \zeta_k = (\sqrt{d} + P_k)/Q_k, \quad \ldots$$

where $a_k Q_k = P_k + P_{k+1}$ and $d - (P_{k+1})^2 = Q_k Q_{k+1}$. Thus the partial quotients satisfy $a_k < 2\sqrt{d}$ and we have the estimate $B_{k+1}/B_k < 1 + 2\sqrt{d}$; such an estimation does not hold in general but it does hold for quadratic surds.

From the upper estimate in Lemma 1.2 we see that $\|y\sqrt{d}\|$ is maximized for integers of the form

$$y = a_{n+1} B_n + a_{n+3} B_{n+2} + \ldots + a_{n+2m+1} B_{n+2m} = B_{n+2m+1} - B_{n-1}.$$

For such numbers

$$\|y\sqrt{d}\| y < |a_{n+1} D_n - D_{n+1}| y = |-D_{n-1}| y < B_{n+2m+1}/B_n < (1 + 2\sqrt{d})^{2m+1}$$

and we have shown

THEOREM 2. *The positive integer solutions of $x^2 - dy^2 = N$ where $2K_1\sqrt{d} < |N| < 2K_2\sqrt{d}$ are given by (2) and the corresponding sums of $A_k$'s with $C_1 \log K_1 < m < C_2 \log(2K_2 + 1) + 3$ where $C_1$ depends only on $d$ and $C_2$ is an absolute constant.*

We conclude with an explicit calculation of the values represented by $x^2 - dy^2$ for $x = A_k + cA_{k+1}$ and $y = B_k + cB_{k+1}$ where $1 \leqslant c \leqslant a_{k+2} - 1$. Since $x^2 - dy^2 = (y\sqrt{d} - x)(y\sqrt{d} - x)^*$ where $(z)^*$ denotes the conjugate of $z$, we have that for $x = A_k$ and $y = B_k$, $x^2 - dy^2 = D_k(D_k)^* = (-1)^{k+1} Q_{k+1}$ (see Perron [2]). Thus for $x = A_k + cA_{k+1}$ and $y = B_k + cB_{k+1}$,

$$x^2 - dy^2 = D_k(D_k)^* (1 - c/\zeta_{k+2})(1 - c/(\zeta_{k+2})^*)$$
$$= (-1)^{k+1} (Q_{k+1} + cQ_{k+2}(2P_{k+2}/Q_{k+2} - c))$$
$$= (-1)^{k+1} V_{k+1}(c).$$

Since $2P_{k+2}/Q_{k+2} = \zeta_{k+2} + (\zeta_{k+2})^*$ we have $\zeta_{k+2} > 2P_{k+2}/Q_{k+2} > \zeta_{k+2} - 1 > c$ and so $V_{k+1}(c) > Q_{k+1}$ for $1 \leqslant c \leqslant a_{k+2} - 1$. The maximum of $V_{k+1}$ occurs when $c = \|P_{k+2}/Q_{k+2}\|$ and is approximately $Q_{k+1} + (P_{k+2})^2/Q_{k+2} = d/Q_{k+2}$. Since $Q_{k+2} = 1$ at the end of each period of $\sqrt{d}$, $V_{k+1}$ can take values as large as $d$.

For an integer of the form (2) with $m > 2$, we see from the recursion formula $B_{k+1} = a_{k+1} B_k + B_{k-1}$ that it can be rewritten as $sB_k + tB_{k+1}$ where $s > 1$ and $t > 0$ are integers. Then the previous calculations may be repeated to find the value of $x^2 - dy^2$ in terms of $Q_{k+1}, Q_{k+2}$ and $P_{k+2}$.

### References

[1] Alexander Ostrowski, *Bemerkungen zur Theorie der Diophantischen Approximationen*, Abh. Math. Sem. Hamburg Univ. 1 (1922), pp. 77–98.

[2] Oskar Perron, *Die Lehre von den Kettenbrüchen*, Bd. I, 3te Aufl., B. G. Teubner, 1954.

[3] Peter Szüsz, *Bemerkungen zur Approximationen einer Reellen Zahl durch Brüche*, Acta Math. Acad. Sci. Hung. 6 (1955), pp. 203–212.

DEPARTMENT OF MATHEMATICS
SUNY AT STONY BROOK
STONY BROOK, NEW YORK 11794 USA

# A generalization of Atkinson's formula to L-functions

by

TOM MEURMAN (Turku)

**1. Introduction.** Let

$$I(q, T) = \sum_{\chi \bmod q} \int_0^T |L(\tfrac{1}{2} + it, \chi)|^2 \, dt \tag{1.1}$$

and define $E(q, T)$ via the identity

$$I(q, T) = \frac{\varphi^2(q)}{q} T \left( \log \frac{qT}{2\pi} + \sum_{p \mid q} \frac{\log p}{p - 1} + 2\gamma - 1 \right) + E(q, T), \tag{1.2}$$

where $\varphi$ is Euler's function and $\gamma$ is his constant.

Consider first the case $q = 1$. Atkinson [1] has established for $E(1, T)$ a very precise explicit expression in terms of two sums involving the divisor function $d(n)$. Recently Jutila [7] found a new interesting application of this formula by showing that it yields in a simple manner Balasubramanian's [2] estimate $E(1, T) \ll T^{1/3 + \varepsilon}$, valid for any positive $\varepsilon$.

The more general function $E(q, T)$ has been studied by Rane [9] (in fact, he considers $E(q, T) - E(q, 1)$) who proved

$$E(q, T) \ll qT^{1/2} \log T. \tag{1.3}$$

A simpler proof of this is due to Balasubramanian and Ramachandra [3].

Our object is to generalize Atkinson's formula to $E(q, T)$ (Theorem 1). Then we deduce by Jutila's method a new inequality for $E(q, T)$ (Corollary 1). In turn, this implies immediately new mean value estimates for $L$-functions (Corollary 2), which can be applied to estimate the density of the zeros in small rectangles (Corollary 3).

We proceed to state the main results. Let

$$e(T, u) = \left( 1 + \frac{\pi u}{2T} \right)^{-1/4} \left( \left( \frac{2T}{\pi u} \right)^{1/2} \operatorname{arsinh} \left( \left( \frac{\pi u}{2T} \right)^{1/2} \right) \right)^{-1}, \tag{1.4}$$

$$f(T, u) = 2T \operatorname{arsinh} \left( \left( \frac{\pi u}{2T} \right)^{1/2} \right) + (\pi^2 u^2 + 2\pi u T)^{1/2} - \frac{\pi}{4}, \tag{1.5}$$

$$g(T, u) = T \log \frac{T}{2\pi u} - T + 2\pi u + \frac{\pi}{4}. \tag{1.6}$$