

References

- [1] E. Bombieri, *Le grand crible dans la théorie analytique des nombres*, Asterisque 18 (1974), Société Mathématique de France, page 22.
 [2] J. H. Conway and A. J. Jones, *Trigonometric diophantine equations (on vanishing sums of roots of unity)*, Acta Arith. 30 (1976), 229–240.
 [3] J. H. Loxton, *On two problems of R. M. Robinson about sums of roots of unity*, ibid. 26 (1974), 159–174.
 [4] H. B. Mann, *On linear relations between roots of unity*, Mathematika 12 (1965), 107–117.
 [5] A. Schinzel, *Reducibility of lacunary polynomials VIII*, Acta Arith. 50 (1988), 91–106.

Received on 7.8.1987

(1743)

Some number-theoretical properties of generalized sum-of-digit functions

by

GERHARD LARCHER (Salzburg) and ROBERT F. TICHY (Wien)

1. Introduction. There has been a great deal of work in investigating the number-theoretical properties of the sum of digits of positive integers in a given number system. In the special case of a q -ary number system ($q \geq 2$) write n in the digit representation

$$(1.1) \quad n = \sum_{i=0}^{\infty} \varepsilon_i q^i$$

with $\varepsilon_i = \varepsilon_i(q, n) \in \{0, \dots, q-1\}$ and $\varepsilon_i = 0$ for $i > [\log n / \log q]$; $[x]$ denotes, as usual, the greatest integer $\leq x$. Then by a famous result of Delange [3]

$$(1.2) \quad \frac{1}{n} \sum_{k=0}^{n-1} s(q, k) = \frac{q-1}{2} \frac{\log n}{\log q} + nF\left(\frac{\log n}{\log q}\right),$$

where $s(q, k) = \sum_{j=0}^{\infty} \varepsilon_j(q, k)$ denotes the sum of q -ary digits and F is a suitable continuous and nowhere differentiable function with period 1. Exact bounds of the error term $F(\log n / \log q)$ have been given by Drazin and Griffiths [4]. A further precise information on the average value of the sum of q -ary digits is given in a recent paper of Foster [6]. In the case $q = 2$ he proved

$$(1.3) \quad -\frac{2}{13} < \frac{2}{n} \sum_{k=0}^{n-1} s(2, k) - \left\lfloor \frac{\log n}{\log 2} \right\rfloor < 1,$$

where both bounds are best possible. A paper of Stolarsky [12] contains a brief survey of the history of such problems and cites many references.

Other authors, especially French mathematicians investigated certain exponential sums, e.g.

$$(1.4) \quad \sum_{k=0}^{n-1} e^{2\pi i h s(q, k)x} \quad (h \text{ integral, } x \text{ irrational})$$

in connection with the uniform distribution of the sequence $(s(q, n)x)_{n=0}^{\infty}$.

These investigations were initiated by Mendès-France [10] and continued in several articles by Coquet, e.g. [2].

In recent time more general digit depending sums and sequences turned out to be of some importance in various fields of applications; for instance the Rudin-Shapiro sequence (cf. Allouche and Mendès-France [1]) with applications in harmonic analysis and in the theory of automata and the Gray code representation with applications in computer science (cf. Sedgewick [11], Flajolet and Ramshaw [5]). The digits $\gamma_i(n)$ in Gray code representation are given by

$$(1.5) \quad \gamma_i(n) = \varepsilon_i(2, n) + \varepsilon_{i+1}(2, n) \bmod 2,$$

and $G(k) = \sum_{j=0}^{\infty} \gamma_j(k)$ denotes the sum of Gray code digits. Obviously $G(k)$ is the number of maximal 0-blocks and 1-blocks in the binary representation of k . In an appendix we prove an explicit formula for

$$(1.6) \quad \frac{1}{n} \sum_{k=0}^{n-1} G(k) - \frac{1}{2} \left\lfloor \frac{\log n}{\log q} \right\rfloor$$

(cf. Foster [6]). From this formula it is possible to derive lower and upper bounds for the expression (1.6) which can be used to give estimates for the average case complexity

$$\frac{n}{2} + 2 \sum_{k \geq 1} G(k) \left(\binom{2n}{n-k} / \binom{2n}{n} \right)$$

of Batcher's sorting algorithm with n files (cf. Sedgewick [11], Flajolet and Ramshaw [5]).

Our main results are concerned with estimates for exponential sums of type (1.4) with respect to Gray code and some extensions. The sequence $(G(n))_{n=0}^{\infty}$ is a special case of a more general class of sequences. Let $q, \kappa \in \mathbb{N}$ (positive integers), $q \geq 2$ and for $i_1, \dots, i_{\kappa} \in \{0, \dots, q-1\}$ let $a(i_1, \dots, i_{\kappa})$ be a real number and assume $a(0, \dots, 0) = 0$. For n in q -ary digit representation (1.1) we define

$$(1.7) \quad t(n) := \sum_{i=0}^{\infty} a(\varepsilon_i, \varepsilon_{i+1}, \dots, \varepsilon_{i+\kappa-1}).$$

In the case $\kappa = 1$ and $a(i) = i$, $t(n) = s(q, n)$ is the usual sum of q -ary digits. For $q = \kappa = 2$ and $a(0, 0) = a(1, 1) = 0$, $a(0, 1) = a(1, 0) = 1$ we have $t(n) = G(n)$. In the case $q = \kappa = 2$ and $a(0, 0) = a(0, 1) = a(1, 0) = 0$, $a(1, 1) = 1$, $t(n)$ is the number of (11)-blocks in the binary representation of n (Rudin-Shapiro-sequence).

For a pleasant formulation of the main theorems we need the following

quantity ω . Let b, l be integers with $0 \leq b, l < q^{\kappa}$ in q -ary representation

$$b = \sum_{i=0}^{\kappa-1} b_i q^i, \quad l = \sum_{i=0}^{\kappa-1} l_i q^i.$$

Then for real x we define

$$(1.8) \quad \sigma = \sigma(x, l, b) = x \cdot \left(t(l) + \sum_{m=1}^{\kappa-1} (a(b_m, \dots, b_{\kappa-1}, l_0, \dots, l_{m-1}) - a(b_m, \dots, b_{\kappa-1}, 0, \dots, 0)) \right),$$

Furthermore we set

$$(1.9) \quad \omega(x) = \max_{b, l} \|\sigma(x, l, b)\| \quad (\|y\| = \min(y - [y], 1 - y + [y])),$$

$$S(N, x) = \sum_{k=0}^{N-1} e^{2\pi i t(k) \cdot x}.$$

THEOREM 1. For all positive integers N and real x

$$|S(N, x)| \leq C_1 N^{\frac{\log(q - C_2 \omega^2(x))}{\log q}}$$

with some positive constants C_1, C_2 only depending on q and κ .

Our second result is concerned with a special lower bound for $S(N, x)$.

THEOREM 2. There is a constant $\delta > 0$ such that for all reals x with $\omega(x) \leq \delta$ and for all $N = q^{j\kappa}$ ($j \in \mathbb{N}$) we have

$$|S(N, x)| \geq C_3 N^{\frac{\log(q - C_4 \omega^2(x))}{\log q}}$$

with some positive constants C_3, C_4 only depending on q and κ .

We can use the above theorems to generalize results concerning the discrepancy of the sequence $(s(q, n)x)_{n=0}^{\infty}$ (cf. [9], [13]) by giving best possible estimates for the discrepancy of the sequence $\tau = (t(n)x)_{n=0}^{\infty}$. Let us recall that the discrepancy of a sequence $\xi = (x_n)_{n=0}^{\infty}$ of real numbers is defined by

$$(1.10) \quad D_N(\xi) = \sup_{0 \leq \alpha < \beta \leq 1} \left| \frac{A(N; \alpha, \beta, x_n)}{N} - (\beta - \alpha) \right|$$

with $A(N; \alpha, \beta, x_n) = \text{card} \{0 \leq n < N: \alpha \leq x_n - [x_n] < \beta\}$. The sequence $\xi = (x_n)_{n=0}^{\infty}$ is called *uniformly distributed modulo 1* if $D_N(\xi)$ tends to 0 (for $N \rightarrow \infty$); cf. the monographs [7], [8].

THEOREM 3. The sequence $\tau = (t(n)x)_{n=0}^{\infty}$ is uniformly distributed mod 1 if

and only if $\omega(hx) \neq 0$ for all positive integers h . If

$$\omega(hx) \geq C_5/h^\eta$$

for all $h \in \mathbb{N}$ and fixed constants $C_5, \eta > 0$, then we have for all $N \in \mathbb{N}$

$$(1.11) \quad D_N(\tau) \leq C_6/(\log N)^{1/2\eta}$$

with a constant C_6 only depending on q, κ, C_5, η . Conversely we have

$$(1.12) \quad D_N(\tau) \geq C_7/(\log N)^{1/2\eta}$$

for infinitely many $N \in \mathbb{N}$ provided that

$$\omega(hx) \leq C_8/h^\eta \quad \text{for infinitely many } h \in \mathbb{N}.$$

Coquet [2] was heavily interested in the distribution behaviour of the sequences $(s(q, n+k)x)_{n=0}^\infty$ uniformly in k . In this special case a (not best-possible) estimate for the uniform discrepancy

$$(1.13) \quad \tilde{D}_N(\xi) = \sup_{k=0,1,2,\dots} D_N(\xi^{(k)}) \quad (\xi^{(k)} = (x_{n+k})_{n=0}^\infty)$$

is known (cf. [14]).

In the following theorem a best possible estimate for $\tilde{D}_N(\tau)$ is established.

THEOREM 4. For all positive integers N and $\tau = (t(n)x)_{n=0}^\infty$ the estimate

$$\tilde{D}_N(\tau) \leq 4 \cdot q^\kappa \max_{1 \leq j \leq N} D_j(\tau)$$

holds.

Remark 1. As an immediate consequence of Theorem 4 we get the estimates (1.11) and (1.12) even for the uniform discrepancy $\tilde{D}_N(\tau)$.

Remark 2. In the case $\kappa = 1$ we have

$$\omega(x) = \max_{i=1,\dots,q-1} \|a(i)x\|;$$

hence the theorem in [9] is a special case of our results.

Remark 3. In the case of Gray code we obtain $\omega(x) = \max(\|x\|, \|2x\|)$. Hence $\chi = (G(h)x)_{n=0}^\infty$ is uniformly distributed mod 1 for all irrationals x ; χ is even well distributed in this case. If x is of approximation type η then we have

$$(1.14) \quad \tilde{D}_N(\tau) \leq \frac{C_9}{(\log N)^{1/2\eta}}$$

with a constant C_9 only depending on η , and this estimate is best possible.

Remark 4. In the case of the Rudin-Shapiro sequence we also obtain $\omega(x) = \max(\|x\|, \|2x\|)$. Hence the same conclusions as in Remark 3 are valid.

2. Proof of Theorem 1. We make use of the following auxiliary results.

LEMMA 1. There are constants $c_1, c_2 > 0$ such that for all $\mu > 0$ and all $z_1, \dots, z_s \in \mathbb{C}$ ($s \geq 2$) we have

$$(2.1) \quad |z_1 + \dots + z_s| \geq (s - c_1 \mu^2) \min_{i=1,\dots,s} |z_i|$$

provided that $\max_{i,j=1,\dots,s} \|\arg z_i - \arg z_j\| \leq \mu$;

$$(2.2) \quad |z_1 + \dots + z_s| \leq (s - c_2 \mu^2) \max_{i=1,\dots,s} |z_i|$$

provided that $\max_{i,j=1,\dots,s} \|\arg z_i - \arg z_j\| \geq \mu$. ($\arg r \cdot e^{2\pi i \gamma} := \gamma$ for $-1/2 < \gamma \leq 1/2$.)

We omit the easy proof.

LEMMA 2. Let $a_1(0), \dots, a_s(0)$ and α_{ik} be complex numbers with $|a_i(0)| = |\alpha_{ik}| = 1$ and $\alpha_{1i} = 1$ ($i, k = 1, \dots, s$; $s \geq 2$). Furthermore assume that

$$(2.3) \quad a_l(j) = \sum_{i=1}^s \alpha_{il} a_i(j-1) \quad \text{for } l = 1, \dots, s \text{ and } j \in \mathbb{N}$$

and set

$$v = \max_{i,j=1,\dots,s} |\arg \alpha_{ik}|.$$

Then there are constants $c_3, c_4 > 0$ (only depending on s) such that

$$(2.4) \quad |a_l(j)| \leq c_3 (s - c_4 v^2)^j \quad \text{for all } l = 1, \dots, s \text{ and } j \in \mathbb{N}.$$

Proof. We proceed by induction. For $j = 0, 1, 2$ the assertion is trivially true and we assume that it is proved for $j \leq n$ ($n \geq 2$). Then by Lemma 1

$$|a_1(n-1)| \leq c_3 (s - c_5 v^2) (s - c_4 v^2)^{n-2}$$

provided that

$$\max_{l,m} (\|\arg a_l(n-2) - \arg a_m(n-2)\|) \geq v/8.$$

Hence

$$(2.5) \quad |a_l(n)| \leq c_3 (s-1) (s - c_4 v^2)^{n-1} + c_3 (s - c_5 v^2) (s - c_4 v^2)^{n-2} \\ \leq c_3 \left(s - \frac{c_5 - c_4}{s} v^2 \right) (s - c_4 v^2)^{n-1}.$$

Now we consider the case

$$\max_{l,m} (|\arg a_l(n-2) - \arg a_m(n-2)|) < \nu/8$$

and without loss of generality we may assume $\arg \alpha_{s1} = \nu$. Then

$$|\arg(\alpha_{si} \cdot a_i(n-2)) - \arg(\alpha_{s1} \cdot a_1(n-2))| \geq \nu/4 - \nu/8 = \nu/8$$

provided that $\arg \alpha_{si} \leq 3\nu/4$ for some i .

Therefore we obtain as above

$$(2.6) \quad |a_s(n)| \leq c_3 \left(s - \frac{c_6 - c_4}{s} \nu^2 \right) (s - c_4 \nu^2)^{n-1}.$$

In the case $\arg \alpha_{si} > 3\nu/4$ (for all $i = 1, \dots, s$) we have

$$|\arg(a_s(n-1)) - \arg(a_1(n-1))| \geq \nu - \nu/4 - 2\nu/8 = \nu/2.$$

Hence we obtain again

$$(2.7) \quad |a_1(n)| \leq c_3 (s - c_7 \nu^2) (s - c_4 \nu^2)^{n-1}.$$

Combining suitable constants c_1, c_8 yields

$$\begin{aligned} |a_1(n+1)| &\leq c_3 (s-1) (s - c_4 \nu^2)^n + c_3 (s - c_8 \nu^2) (s - c_4 \nu^2)^{n-1} \\ &\leq c_3 \left(s - \frac{c_8 - c_4}{s} \nu^2 \right) (s - c_4 \nu^2)^n \leq c_3 (s - c_4 \nu^2)^{n+1}. \end{aligned}$$

Thus the proof of Lemma 2 is complete.

Now we continue with the proof of Theorem 1. Let N be a positive integer with q^x -ary digit representation

$$(2.8) \quad N = \sum_{j=0}^r N_j q^{xj}, \quad N_j \in \{0, 1, \dots, q^x - 1\}, \quad N_r \neq 0.$$

We have (using the notation $L_j = N_{j+1} q^{x(j+1)} + \dots + N_r q^{xr}$ for $j = 0, \dots, r-1$ and $L_r = 0$)

$$\begin{aligned} (2.9) \quad S(N, x) &= \sum_{j=0}^r \sum_{\varepsilon=0}^{N_j-1} \sum_{k=L_j+\varepsilon q^{xj}}^{L_j+(\varepsilon+1)q^{xj}-1} e^{2\pi i t(k)x} \\ &= \sum_{j=0}^r \sum_{\varepsilon=0}^{N_j-1} \sum_{l=0}^{q^x-1} \sum_{k=lq^{x(j-1)}}^{(l+1)q^{x(j-1)}-1} \exp_l(2\pi i t(k+L_j+\varepsilon q^{xj})x) \end{aligned}$$

($\exp(t) = e^t$).

Let $\varepsilon = \sum_{j=0}^{x-1} \varepsilon_j q^j$, $l = \sum_{j=0}^{x-1} l_j q^j$ be given in q -ary digit representation and

let k be given such that $l \cdot q^{x(j-1)} \leq k < (l+1)q^{x(j-1)}$, then

$$(2.10) \quad t(k+L_j+\varepsilon q^{xj}) = t(k) + t(L_j+\varepsilon q^{xj}) + \sum_{m=1}^{x-1} (a(l_m, \dots, l_{x-1}, \varepsilon_0, \dots, \varepsilon_{m-1}) - a(l_m, \dots, l_{x-1}, 0, \dots, 0)).$$

We use the notation

$$\varphi(l, \varepsilon, j, x) = \exp((t(k+L_j+\varepsilon q^{xj}) - t(k))x),$$

and derive from (2.9)

$$(2.11) \quad S(N, x) = \sum_{j,\varepsilon,l} \varphi(l, \varepsilon, j, x) \sum_{k=lq^{x(j-1)}}^{(l+1)q^{x(j-1)}-1} e^{2\pi i t(k)x}.$$

Now we consider the sum

$$S_l(j, x) = \sum_{k=lq^{xj}}^{(l+1)q^{xj}-1} e^{2\pi i t(k)x} \quad \text{for } 0 \leq l \leq q^x - 1$$

and obtain

$$\begin{aligned} (2.12) \quad S_l(j, x) &= \sum_{b=0}^{q^x-1} \sum_{k=bq^{x(j-1)}}^{(b+1)q^{x(j-1)}-1} e^{2\pi i t(k+lq^{xj})x} \\ &= \sum_{b=0}^{q^x-1} \sigma(x, l, b) S_b(j-1, x) \end{aligned}$$

(in the notation (1.8)). By Lemma 2 we have for $0 \leq l \leq q^x - 1$

$$|S_l(j, x)| \leq c_3 (q^x - c_4 \omega^2(x))^j,$$

and therefore

$$\begin{aligned} |S(N, x)| &\leq c_9 \sum_{j=0}^r (q^x - c_4 \omega^2(x))^{j-1} \\ &= c_9 \frac{(q^x - c_4 \omega^2(x))^r - 1}{q^x - c_4 \omega^2(x) - 1} \leq c_{10} (q^x - c_4 \omega^2(x))^{rx} \\ &\leq c_{10} \left(q - \frac{c_4}{x q^x} \omega^2(x) \right)^{xr} = c_1 N^{\frac{\log(q - c_4 \omega^2(x))}{\log q}}. \end{aligned}$$

Thus the proof of Theorem 1 is complete.

3. Proof of Theorem 2. We make use of the following auxiliary results.

LEMMA 3. For $z_j = r_j e^{2\pi i \alpha_j}$, $q_j = c^{2\pi i \beta_j}$ ($j = 1, \dots, s$; $s \geq 2$) with $|\alpha_j|,$

$|\beta_j| \leq \pi/8$ let $B = z_1 + \dots + z_s$ and $B_q = \varrho_1 z_1 + \dots + \varrho_s z_s$. Then

$$|\arg B_q - \arg B| \leq 2 \max_{j=1, \dots, s} |\beta_j|.$$

Proof. Without loss of generality we assume

$$\tan(\arg B) = \frac{\sum_{i=1}^s r_i \sin \alpha_i}{\sum_{i=1}^s r_i \cos \alpha_i}, \quad \tan(\arg B_q) = \frac{\sum_{i=1}^s r_i \sin(\alpha_i + \beta_i)}{\sum_{i=1}^s r_i \cos(\alpha_i + \beta_i)}$$

and

$$\arg B_q - \arg B \leq \tan(\arg B_q) - \tan(\arg B) = \frac{\sum_{i,j=1}^s r_i r_j \sin(\alpha_i - \alpha_j + \beta_j)}{\sum_{i,j=1}^s r_i r_j \cos \alpha_i \cos(\alpha_j + \beta_j)}.$$

Because of

$$|\alpha_i - \alpha_j| \leq \pi/4, \quad |\beta_j| \leq \pi/4 \quad \text{and} \quad |\alpha_j + \beta_j| \leq \pi/4$$

we have

$$\sin(\alpha_i - \alpha_j + \beta_j) \leq \sin(\alpha_i - \alpha_j) + |\beta_j|$$

and

$$\cos \alpha_i \cos(\alpha_j + \beta_j) \geq 1/2.$$

Hence

$$\begin{aligned} \arg B_q - \arg B &\leq \frac{\sum_{i,j=1}^s r_i r_j \sin(\alpha_i - \alpha_j)}{\sum_{i,j=1}^s r_i r_j \cos \alpha_i \cos(\alpha_j + \beta_j)} + 2 \frac{\sum_{i,j=1}^s r_i r_j |\beta_j|}{\sum_{i,j=1}^s r_i r_j} \\ &\leq 2 \max_{i=1, \dots, s} |\beta_i|, \end{aligned}$$

and the proof of Lemma 3 is complete.

LEMMA 4. There are constants $\delta, c_1, c_2 > 0$ depending only on s such that whenever $|\arg a_i(0)| < \delta$ and $|\arg \alpha_{ik}| < \delta$ for all i, k then for all l and j

$$(3.1) \quad |a_l(j)| \geq c_1 (s - c_2 v^2)^j$$

(in the notation of Lemma 2).

Proof. If δ is small enough the assertion is trivially true for $j = 0, 1$. We proceed by induction and assume that (3.1) is true for all $j \leq n$ ($n \geq 1$).

We obtain by Lemma 3 for sufficiently small δ

$$\|\arg(a_l(m)) - \arg(a_1(m))\| \leq 2v$$

for all l and $m = 1, 2, 3, \dots$. Hence by the induction hypothesis and Lemma 1 we have

$$|a_l(n+1)| \geq c_1 (s - c_2 v^2) (s - c_2 v^2)^n = c_1 (s - c_2 v^2)^{n+1}$$

which proves Lemma 4.

By the recurrence (2.12) and Lemma 4 we obtain for all reals x with $\omega(x) < \delta$ and all $N = q^{xj}$ that

$$|S(N, x)| \geq C_3 N^{\frac{\log(q - C_4 \omega^2(x))}{\log q}}$$

which completes the proof of Theorem 2.

4. Proof of Theorems 3 and 4. The first assertion of Theorem 3 follows immediately from Theorems 1, 2 and Weyl's criterion for uniform distribution of sequences (cf. [10]). The upper bound (1.11) can be established by means of the Erdős-Turán inequality

$$(4.1) \quad D_N(\xi) \leq 6 \left(\frac{1}{H} + \sum_{h=1}^H \frac{1}{h} \left| \frac{1}{N} \sum_{n=0}^{N-1} e^{2\pi i h x_n} \right| \right) \quad (\xi = (x_n))$$

(choosing a suitable positive integer H) by verbally the same calculations as in [9] and [13]. The lower bound (1.12) follows from Koksma's inequality

$$(4.2) \quad D_N(\xi) \geq \frac{1}{2\pi h N} \left| \sum_{n=0}^{N-1} e^{2\pi i h x_n} \right|$$

(cf. [8]) as in the special cases [9], [13].

Finally we establish a proof of Theorem 4. For $N, k \in \mathbb{N}$ let a, r be non-negative integers such that $q^{x(r-1)} \leq N < q^{xr}$ and $aq^{xr} \leq k < (a+1)q^{xr}$. First we consider the case

$$(4.3) \quad aq^{xr} \leq k < (k+N-1) < (a+1)q^{xr}.$$

Let

$$b_0 q^{x(r-1)} + aq^{xr} \leq k < (b_0+1)q^{x(r-1)} + aq^{xr}$$

and

$$(b_0 + B_0 q^{x(r-1)} + aq^{xr} \leq k + N - 1 < (b_0 + B_0 + 1)q^{x(r-1)} + aq^{xr}$$

with $0 \leq b_0 < b_0 + B_0 < q^x$. Then the sequence $(t(h+k)x)_{n=0}^{N-1}$ consists of the following parts

$$(t(n+k)x)_{n=0}^{(b_0+1)q^{x(r-1)} + aq^{xr} - k - 1},$$

$$(4.4) \quad (t(n+k)x)_{n=b_0q^{x(r-1)}+aq^{xr}-k}^{(b+1)q^{x(r-1)}+aq^{xr}-k-1} \quad \text{for } b = b_0+1, \dots, b_0+B_0-1$$

and

$$(t(n+k)x)_{n=(b_0+B_0)q^{x(r-1)}+aq^{xr}-k}^{N-1}.$$

Each of these subsequences is of the form $\xi = (\alpha + t(n)x)_{n=n_0}^{n_1}$ where $0 \leq \alpha < 1$ and $0 < n_0 \leq n_1 \leq q^{x(r-1)}$. Denoting by $D_{n_0, n_1}(\xi)$ the discrepancy of the sequence ξ we obtain by well known (and simple) properties of the discrepancy

$$(n_1 - n_0 + 1)D_{n_0, n_1}(\xi) \leq n_1 D_{n_1}(\tau) + (n_0 - n_1)D_{n_0}(\tau).$$

Hence by [8, Theorem 2.6, p. 115] we have

$$ND_N(\tau^{(k)}) \leq 2q^x q^{x(r-1)} \max_{0 \leq j \leq q^{x(r-1)}} D_j(\tau);$$

thus

$$(4.5) \quad \tilde{D}_N(\tau) \leq 2q^x \max_{0 \leq j \leq N} D_j(\tau).$$

In the case $aq^{xr} \leq k < (a+1)q^{xr} \leq k+N-1$ we divide the sequence $(t(n+k)x)_{n=0}^{N-1}$ into two parts

$$(4.6) \quad (t(n+k)x)_{n=0}^{(a+1)q^{xr}-k-1} \quad \text{and} \quad (t(n+(a+1)q^{xr}))_{n=0}^{k+N-(a+1)q^{xr}-1}.$$

Both sequences satisfy the condition of the above case and therefore

$$(4.7) \quad \tilde{D}_N(\tau) \leq 4q^x \max_{0 \leq j \leq N} D_j(\tau).$$

Thus the proof of Theorem 4 is complete.

5. Appendix. In the following we establish an "explicit" formula for

$$S(n) = \sum_{k=0}^{n-1} G(k) - \frac{n}{2} \left[\frac{\log n}{\log k} \right].$$

First we note that

$$(5.1) \quad G(2^k + n) = \begin{cases} G(n) + 2 & \text{for } 0 \leq n < 2^{k-1}, \\ G(n) & \text{for } 2^{k-1} \leq n < 2^k. \end{cases}$$

Applying (5.1) we obtain for

$$A(n) = \sum_{k=1}^{n-1} G(k)$$

the identity

$$\begin{aligned} A(2^s) &= \sum_{1 \leq r < 2^{s-1}} G(r) + \sum_{2^{s-1} \leq r < 2^s} G(r) \\ &= A(2^{s-1}) + \sum_{2^{s-1} \leq r < 2^{s-2} + 2^{s-1}} G(r) + \sum_{2^{s-1} + 2^{s-2} \leq r < 2^s} G(r) \\ &= A(2^{s-1}) + \sum_{0 \leq r < 2^{s-1}} G(r) + 2 \sum_{0 \leq r < 2^{s-1}} 1 \\ &= 2A(2^{s-1}) + 2^{s-1}. \end{aligned}$$

Hence, by induction,

$$(5.2) \quad A(2^s) = \frac{s}{2} 2^s.$$

Every odd number has a unique representation of the form

$$(5.3) \quad n_m = 2^{t_0} + 2^{t_0+t_1} + \dots + 2^{t_0+\dots+t_m}$$

with integers $m \geq 0$, $t_0 = 0$, $t_j \geq 1$ ($j \geq 1$). Furthermore we set

$$(5.4) \quad n_0 = 1, \quad n_i = 1 + 2^{t_1} + \dots + 2^{t_1+\dots+t_i} \quad \text{for } 1 \leq i \leq m.$$

Then we have

$$\begin{aligned} (5.5) \quad A(n_m) &= A(2^{t_1+\dots+t_m}) + \sum_{2^{t_1+\dots+t_m} \leq r < n_m} G(r) \\ &= A(2^{t_1+\dots+t_m}) + \sum_{0 \leq r < n_{m-1}} G(2^{t_1+\dots+t_m} + r). \end{aligned}$$

In the case $t_m = 1$ (i.e. $n_{m-1} > 2^{t_1+\dots+t_{m-1}}$) we obtain by (5.1) and (5.5) ($m \geq 2$)

$$\begin{aligned} (5.6) \quad A(n_m) &= A(2^{t_1+\dots+t_m}) + A(n_{m-1}) + 2 \sum_{0 \leq r < 2^{t_1+\dots+t_{m-1}}} 1 \\ &= A(2^{t_1+\dots+t_m}) + A(n_{m-1}) + 2(n_{m-1} - n_{m-2}). \end{aligned}$$

In the case $t_m > 1$ (i.e. $n_{m-1} < 2^{t_1+\dots+t_{m-1}}$) we obtain by (5.1) and (5.5)

$$\begin{aligned} (5.7) \quad A(n_m) &= A(2^{t_1+\dots+t_m}) + \sum_{0 \leq r < n_{m-1}} (G(r) + 2) \\ &= A(2^{t_1+\dots+t_m}) + A(n_{m-1}) + 2n_{m-1}. \end{aligned}$$

Combining (5.6) and (5.7) gives

$$(5.8) \quad A(n_m) = A(2^{t_1+\dots+t_m}) + A(n_{m-1}) + 2n_{m-1} - 2\delta_m n_{m-2} \quad (m \geq 1),$$

where

$$\delta_m = \begin{cases} 0 & \text{for } t_m > 1, \\ 1 & \text{for } t_m = 1, \end{cases}$$

$$n_{-1} = 0.$$

Summing up and applying (5.2) (note that $A(n_0) = 0$) yields

$$\begin{aligned} A(n_m) &= \sum_{r=1}^m (A(2^{t_1+\dots+t_r}) + 2n_{r-1}) - 2 \sum_{r=1}^m \delta_r n_{r-2} \\ &= \sum_{r=1}^m \left(\frac{t_1 + \dots + t_r}{2} 2^{t_1+\dots+t_r} + 2n_{r-1} \right) - 2 \sum_{r=1}^m \delta_r n_{r-2} \\ &= \sum_{r=1}^m \left(\frac{t_1 + \dots + t_r}{2} (n_r - n_{r-1}) + 2n_{r-1} \right) - 2 \sum_{r=1}^m \delta_r n_{r-2} \\ &= \sum_{r=1}^m \left(2 - \frac{t_r}{2} \right) n_{r-1} + \frac{(t_1 + \dots + t_m) n_m}{2} - 2 \sum_{r=1}^m \delta_r n_{r-2}. \end{aligned}$$

Observing that $t_r + \dots + t_m = \left\lceil \frac{\log n_m}{\log 2} \right\rceil$ we have proved the following formula

$$(5.10) \quad \frac{S(n_m)}{n_m} = \frac{1}{n_m} \sum_{r=1}^m \left(\left(2 - \frac{t_r}{2} \right) n_{r-1} - 2\delta_r n_{r-2} \right).$$

It should be remarked that the expression (5.10) does not change its value if n_m is replaced by $2^\beta n_m$ (β an arbitrary positive integer). Hence the explicit formula (2.10) is valid also in the general case (not only for odd numbers n_m).

Since $n_{r-1}/n_r \leq 1/2$ (for every $r \geq 1$), an application of Foster's lower bound [6, Theorem 1] yields

$$(5.11) \quad \frac{S(n_m)}{n_m} = \frac{1}{2n_m} \sum_{r=1}^m (2 - t_r) n_{r-1} + \frac{1}{n_m} \sum_{r=1}^m (n_{r-1} - 2\delta_r n_{r-2}) \geq -\frac{1}{13}.$$

An elementary observation shows that $S(n_m)/n_m$ increases if all $t_r > 2$ are replaced by $t_r = 2$. Hence, for determining an upper bound for $S(n_m)/n_m$ it suffices to consider the case $t_r = 1$ or 2 ($r = 1, \dots, m$). We prove $S(n_m)/n_m \leq 7/10$ by induction; the cases $m = 0, 1$ are trivial. Assuming this bound for $m-1$ we obtain in the case $t_m = 2$ by (5.10)

$$\begin{aligned} (5.12) \quad \frac{S(n_m)}{n_m} &\leq \frac{1}{n_m} \left(S(n_{m-1}) + \left(2 - \frac{t_m}{2} \right) n_{m-1} - 2\delta_m n_{m-2} \right) \\ &\leq \frac{7}{10} \frac{n_{m-1}}{n_m} + \frac{n_{m-1}}{n_m} \leq \frac{7}{10}, \end{aligned}$$

since $n_{m-1}/n_m \leq 1/(1 + 2^{t_m-1})$.

In the case $t_m = 1$ we have

$$\frac{S(n_m)}{n_m} = \frac{n_{m-1}}{n_m} \left(\frac{S(n_{m-1})}{n_{m-1}} - \frac{1}{2} \right) + 2 \frac{n_{m-1} - n_{m-2}}{n_m}.$$

The expression $(n_{m-1} - n_{m-2})/n_m$ takes its maximal value if $t_1 = \dots = t_l = 1$ and $t_{l+1} = \dots = t_{m-1}$ (for some l with $1 \leq l \leq m-1$). Hence after a simple calculation

$$\frac{n_{m-1} - n_{m-2}}{n_m} = \frac{2^{2m-l-2}}{(2^{l+1} - 1) + \frac{1}{3} 2^{l+2} (4^{m-l-1} - 1) + 2^{2m-l-1}} \leq \frac{3}{10}.$$

From this estimate and the induction hypothesis we derive

$$(5.13) \quad \frac{S(n_m)}{n_m} \leq \frac{1}{2} \left(\frac{7}{10} - \frac{1}{2} \right) + \frac{3}{5} = \frac{7}{10},$$

which completes the induction argument. Combining (5.11) and (5.13) we have for all positive integers n (note the remark after (5.10))

$$(5.14) \quad -\frac{1}{13} \leq \frac{S(n)}{n} \leq \frac{7}{10}.$$

These bounds seem to be quite far away from the optimal bounds. It may be possible to derive the optimal bounds from (5.10) using inductive arguments similar to Foster's [6]. However, some numerical calculations become rather extensive in this case.

References

- [1] J. P. Allouche and M. Mendès-France, *On an extremal property of the Rudin-Shapiro sequence*, *Mathematika* 32 (1985), 33–38.
- [2] J. Coquet, *Sur certaines suites uniformément équireparties modulo 1*, *Acta Arith.* 36 (1980), 157–162.
- [3] H. Delange, *Sur la fonction sommatoire de la fonction 'somme des chiffres'*, *Enseignement Math.* 2.21 (1975), 31–47.
- [4] M. P. Drazin and J. S. Griffiths, *On the decimal representation of integers*, *Proc. Cam. Phil. Soc.*, (4), 48 (1952), 555–565.
- [5] P. Flajolet and L. Ramshaw, *Gray code and odd-even merge*, *SIAM J. Comput.* 9 (1980), 142–158.
- [6] D. M. E. Foster, *Estimates for a remainder term associated with the sum of digits function*, to appear in *Glasgow Math. J.*
- [7] E. Hlawka, *Theorie der Gleichverteilung*, *Bibl. Inst. Mannheim-Wien-Zürich*, 1979.
- [8] L. Kuipers and H. Niederreiter, *Uniform distribution of sequences*, *Wiley and Sons*, New York 1974.
- [9] G. Larcher, *Exponential sums of digit-depending sequences and uniform distribution*, to appear.
- [10] M. Mendès-France, *Nombres normaux. Applications aux fonctions pseudo-aléatoires*, *J. Analyse Math.* 20 (1967), 1–56.

- [11] R. Sedgewick, *Data movement in odd-even merging*, SIAM J. Comput. 7 (1978), 239–272.
- [12] K. B. Stolarsky, *Power and exponential sums related to binomial digit parity*, SIAM J. Appl. Math. 32 (1977), 717–730.
- [13] R. F. Tichy and G. Turnwald, *On the discrepancy of some special sequences*, J. Number Theory 26 (1987), 68–78.
- [14] —, — *Gleichmässige Diskrepanzabschätzung für Ziffernsummen*, Anz. Österr. Akad. Wiss. (1986), 17–21.

MATHEMATISCHES INSTITUT
DER UNIVERSITÄT SALZBURG
Hellbrunner Strasse 34
A-5020 Salzburg, Austria

ABTEILUNG FÜR TECHNISCHE MATHEMATIK
TU WIEN
Wiedner Hauptstrasse 8–10
A-1040 Wien, Austria

Received on 14.8.1987

(1744)

On the number of values taken by a polynomial over a finite field

by

J. F. VOLOCH (Rio de Janeiro)

Let F_q be the finite field with q elements and $f(x) \in F_q[x]$ a polynomial of degree n . Let $r(f) = \#f(F_q)$, considering f as a function $f: F_q \rightarrow F_q$. A classical problem, raised by Chowla [3] (see [4] for other references), is to estimate $r(f)$ in terms of n and q . One has the trivial bounds $q/n \leq r(f) \leq q$. The lower bound is essentially best possible and a characterization of the cases with equality when q is prime was obtained in [2].

On the other hand, if f is a “general” polynomial (in a sense that can be made precise, see below) Uchiyama [6] proved that $r(f) \geq q/2 + O(q^{1/2})$ and Birch and Swinnerton-Dyer [1] found the precise result

$$r(f) = q \left(\sum_{i=1}^n \frac{(-1)^{i-1}}{i!} \right) + O(q^{1/2}).$$

They proved this when the Galois group of $f(x) = y$ over $\bar{F}_q(y)$ is the full symmetric group. Of course these results are interesting only when q is large compared to n . The purpose of this paper is to give lower bounds for $r(f)$, valid for f “general”, which improves on the above bounds in several cases.

Uchiyama’s condition is that the polynomial

$$f^*(u, v) = (f(u) - f(v))/(u - v)$$

is absolutely irreducible. When this is the case he could apply Weil’s estimate ([7]) on the number of points of $f^*(u, v) = 0$ over F_q to get his result.

To relate the number of solutions of $f^*(u, v) = 0$ in F_q^2 with $r(f)$, Uchiyama [6] proved the following:

LEMMA 1. *Let N be the number of solutions of $f^*(u, v) = 0$ in F_q^2 and n_0 the number of solutions of $f'(x) = 0$ in F_q . Then*

$$r(f) \geq q^2/(N + q - n_0).$$

Proof. First notice that $f^*(u, v) = 0$ and $u \neq v$ if and only if $f(u) = f(v)$ and that $f^*(u, u) = f'(u)$. Let $\{a_1, \dots, a_r\} = f(F_q)$, $r = r(f)$ and n_i