

COMPLEXITY OF ELIAS ALGORITHM BASED ON CODES
WITH COVERING RADIUS THREE

L. H. ASLANYAN *¹, H. E. DANOYAN **²

¹ Institute for Informatics and Automation Problems of NAS of Armenia

² Chair of Discrete Mathematics and Theoretical Informatics YSU, Armenia

The algorithm for finding the set of "nearest neighbors" in a set using compact blocks and hash functions is known (Elias algorithm). In this paper hash coding schemas associated to coverings by spheres of the same radius are considered. In general, such coverings can be obtained via perfect codes, and some other generalizations of perfect codes such as uniformly packed or quasi perfect codes. We consider the mentioned algorithm for Golay code and for two-error-correcting primitive BCH codes of length $2^m - 1$ for odd m . A formula of time complexity of the algorithm is obtained in these cases.

MSC2010: 68P10; 68P30.

Keywords: nearest neighbors, best match, Elias algorithm, hash functions, quasi-perfect codes, uniformly packed codes, coset weight distribution.

1. Introduction. Let $E = \{0, 1\}$. Consider the Cartesian power E^n , which is known as the set of vertices of n -dimensional unit cube. For $x, y \in E^n$ denote by $d(x, y)$ the Hamming distance between the vectors x and y . For an $x \in E^n$ denote by $S_r^n(x)$ the sphere of radius r centered at point x , i.e. $S_r^n(x) = \{y \in E^n / d(x, y) \leq r\}$ and denote by $O_r^n(x)$ the shell of radius r , i.e. $O_r^n(x) = \{y \in E^n / d(x, y) = r\}$. We will denote by $car(x)$ the carrier of the vector $x = (x_1, \dots, x_n)$, i.e. $car(x) = \{i / x_i \neq 0, i \in \{1, \dots, n\}\}$. Denote by $w(x)$ the weight of the vector x , i.e. $w(x) = \sum_{i=1}^n x_i$. We will call a code a nonempty subset C of E^n (usually some other additional properties take place such as linearity, cyclicity, etc.). The code C is called linear, if C is a linear subspace of E^n . Due to the binary nature of the spaces considered, C is linear when $\forall c_1, c_2 \in C \Rightarrow c_1 + c_2 \in C$. Denote by d_C the minimum distance of code C , i.e.

$$d_C = \min_{c_1, c_2 \in C, c_1 \neq c_2} d(c_1, c_2).$$

The packing radius [1] of code C is called the following number: $d_C = \lfloor \frac{d_C - 1}{2} \rfloor$. Denote by R_C the covering radius of the code C , i.e. $R_C = \max_{x \in E^n} \min_{c \in C} d(x, c)$. In the sequel, when it will make no confusion, we will use notations d, r and R instead of d_C, r_C and R_C respectively. We say that we have an $(n, M, d)R$ code C , if it has length n , cardinality M , distance d and covering radius R . When it is known that C is linear to fix that we use the notation $[n, k, d]R$,

* E-mail: lasl@scil.am

** E-mail: hdanoyan@yandex.ru

where k is the dimension of C as linear subspace. Recall that the code C is called perfect [1], if $r_C = R_C$. It is known [1, 2], that in binary spaces nontrivial perfect codes can have only the following two parameter sets.

- I. $(2^m - 1, 2^{2^m - m - 1}, 3)1$,
- II. $(23, 2^{11}, 7)3$,

where I corresponds to parameters of Hamming code, and II corresponds to parameters of Golay code. For $x \in E^n$ the coset of linear code C is called the set $x + C = \{x + c / c \in C\}$. As it is known [1], two different cosets do not intersect and their union cover the space E^n . We denote by G_C the generator matrix of the linear code C . Recall that G_C is a matrix with rows forming basis of C . Let us denote by H_C the parity check matrix of linear code C . If C is $[n, k, d]R$ code, then H_C is $(n - k) \times k$ matrix for which the equation $c \in C \leftrightarrow H_C c^T = 0$ takes place. For $x \in E^n$ denote by $A_i(x)$ the number of codewords of C located at distance i from x . The nonnegative integers $A_0^C, A_1^C, \dots, A_n^C$, where $A_i^C = |\{c \in C / w(c) = i\}|$ are called weight spectra of code C . Let us denote by $W_C(x)$ the weight enumerator of code C : $W_C(x) = \sum_{i=0}^n A_i^C x^i$. A code C will be called uniformly packed [3], if there are numbers a_1, a_2, \dots, a_{R_C} such that for all $x \in E^n$ the equation $\sum_{i=0}^{R_C} a_i A_i^C(x) = 1$ takes place. Denote by $K_j^n(x)$ the Kravchouk polynomial of degree j [1, 4], i.e.

$$K_j^n(x) = \sum_{i=0}^j (-1)^i \binom{n-x}{i-j} \binom{x}{i}, \text{ where } \binom{x}{j} = \frac{x(x-1)\dots(x-j+1)}{j!}.$$

Denote $L_C(x) = \sum_{i=0}^{R_C} a_i K_i^n(x)$. A code C will be called quasi perfect [1, 4], if $R_C = r_C + 1$. Many families of quasi perfect codes are known for the covering radius ≤ 4 [5–10], but the general problem of existence of quasi-perfect codes by the given parameters is not completely solved yet [5]. When the geometrical interpretation of spherical covers is considered in the models of search of similarities, besides the perfect codes their other possible extensions can be considered and applied, such as quasi perfect codes or uniformly packed codes. The paper is organized as follows: in section 2 is brought definitions and coset weight structures of two error correcting primitive BCH codes and Golay codes, keeping in mind the fact that these are uniformly packed codes. Then in section 3 we consider the Elias algorithm for hash function obtained via these codes and get the formula representation of complexity of the algorithm using coset weight structure of mentioned codes.

2. Preliminaries.

2.1. Coset Weight Distribution of Uniformly Packed Codes. For a linear code C we introduced the coset as the shift of the code. Later we need the coset weight distributions of two error-correcting BCH codes for length $n = 2^{2s+1} - 1$ and for the Golay code. As these codes can be considered as uniformly packed codes [3] we can find mentioned distributions by the method which brought in [3]:

Theorem 1. Let C be uniformly packed code with parameters a_0, a_1, \dots, a_n . Then the polynomial $L_C(x)$ has R distinct roots between 0 and n [3].

Let us denote those roots by ξ_1, \dots, ξ_R . Mention that if C is uniformly packed code containing zero vector, then there exists a uniformly packed code with the same parameters and with the minimum weight b , where $0 \leq b \leq R$ which we denote by C_b . From the proof of the Theorem 1 follows:

Theorem 2. For the weight function of the uniformly packed code C_b takes place the following equality [3]

$$W_{C_b}(x) = \frac{(1+x)^n}{\sum_{i=0}^R a_i \binom{n}{i}} + \sum_{i=1}^R B_{\xi_i}^b (1+x)^{n-\xi_i} (1+x)^{\xi_i}. \quad (1)$$

In (1) $B_{\xi_i}^b$'s are constants, which can be calculated from (1) by equalizing the corresponding coefficients in left and right sides and assuming that we know first R coefficients of $W_{C_b}(x)$. In other words, to find the coefficients $B_{\xi_i}^b$'s we must solve the corresponding linear system of R equations with R variables assuming that we know first R values of weight spectra of the code C_b . From the Theorem 2 follows:

$$A_i^{C_b} = \frac{\binom{n}{i}}{\sum_{j=1}^R a_j \binom{n}{j}} + \sum_{j=1}^R B_{\xi_j}^b K_i^n(x_j). \quad (2)$$

Consequently to know the coset weight distributions of the uniformly packed code, we must calculate only first R coset weights, which are $A_0^{C_b}, A_1^{C_b}, \dots, A_{R-1}^{C_b}$.

2.2 Two-Error-Correcting Primitive BCH Codes. Let us denote the finite field of q elements (q is a power of a prime number) by F_q . We will consider finite fields with characteristic 2 [1]. Denote by α the primitive element of the field F_q . Consider the set of formal polynomials $F_q[x]$ with coefficients from the field F_q . As it is known [1], the factor ring $R[x] = F_q[x]/(x^n - 1)$ is a ring of principal ideals, i.e. each ideal in $R[x]$ is principal. A $[n, k, d]$ BCH code C will be called cyclic, if it is linear and from $c = (c_1, c_2, \dots, c_n) \in C$ follows that $(c_n, c_1, \dots, c_{n-1}) \in C$. We can correspond to each vector (c_1, c_2, \dots, c_n) the polynomial $c_1 + c_2x + \dots + c_nx^{n-1}$, so we can consider a code as a subset of $R[x]$. It is known [1], that each cyclic code is an ideal of $R[x]$ and consequently there is a unique monic polynomial $g(x)$ (generator polynomial) of minimum degree such that $\forall c \in C \exists f(x), c(x) = g(x)f(x)$, where the multiplication is taken in $R[x]$. Two-error-correcting primitive BCH codes are defined as cyclic codes for lengths $n = 2^m - 1$ [1, 4]. The generator polynomial is $g(x) = \text{scm}\{M_{\alpha}(x), M_{\alpha^3}(x)\}$, where by $M_{\alpha^i}(x)$ is denoted the minimal polynomial of the element α^i . It is known, that these codes have $2^m - 2m - 1$ and minimum distance 5 [1]. Also it is known that two-error-correcting primitive BCH codes are quasi-perfect codes [1, 11]. Weight distribution of these codes is calculated in [1, 12]. For odd m two-error-correcting BCH codes are uniformly packed [3] with parameters $a_0 = a_1 = 1, a_2 = a_3 = \frac{6}{n-1}$. Roots of $L_{BC}(x)$ are $\xi_1 = \frac{n+1}{2} - \sqrt{\frac{n+1}{2}}, \xi_2 = \sqrt{\frac{n+1}{2}}$ and $\xi_3 = \frac{n+1}{2} + \sqrt{\frac{n+1}{2}}$. It is known that for odd m there are four distinct weight distributions [3] and for even m there are eight distinct weight distributions brought in [13].

2.3 Golay Code. First let us define the extended Golay code which has length 24. Let A_{11} be the Hadamard matrix of Paley type [1], i.e.

$$A_{11} = \begin{pmatrix} 1 & 1 & 0 & 1 & 1 & 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 & 1 & 1 & 1 & 0 & 0 & 0 & 1 \\ 1 & 0 & 1 & 1 & 0 & 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 1 & 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 1 & 0 & 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 & 1 & 1 & 0 & 1 & 1 & 1 \\ 1 & 0 & 0 & 0 & 1 & 0 & 1 & 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 & 0 & 1 & 0 & 1 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 & 0 & 0 & 1 & 0 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 & 0 & 0 & 0 & 1 & 0 & 1 & 1 \\ 1 & 0 & 1 & 1 & 1 & 0 & 0 & 0 & 1 & 0 & 1 \end{pmatrix}.$$

Let I_{11} be the identity matrix of size 11×11 . The generator matrix of extended Golay code [1, 4] is the following

$$G_{\Gamma_{24}} = \begin{pmatrix} 1_{11}^T & I_{11} & 1_{11}^T & A_{11} \\ 0 & 0_{11} & 1 & 1_{11} \end{pmatrix},$$

where by 0_m and 1_m is denoted the all zero and one vectors respectively of length m . The Golay code Γ_{23} is obtained by deleting the last coordinate from each codeword of extended Golay code. It is known [1, 4], that Γ_{23} is a perfect three-error-correcting [23, 12, 7]₃ code therefore we can consider it as uniformly packed code with parameters $a_0 = a_1 = a_2 = a_3 = 1$. The roots of $L_{\Gamma_{23}}(x)$ are $\xi_1 = 8$, $\xi_2 = 12$ and $\xi_3 = 16$.

3. Elias Algorithm. Let we have a subset (or a file) $F \in E^n$ and a query element $x \in E^n$. Let us consider the problem of finding the set of all "nearest neighbors" of F to x . More precisely it is required to find the set $b(x, F) = \{y \in F / d(x, y) = c\}$, where $c = \min_{z \in F} d(x, z)$. To propose an algorithm for solving the problem of nearest neighbors in applied level, hash coding schemas are considered [14, 15]. We will brought a brief description of such schemes. Hash function is defined as a function $h : E^n \rightarrow V$, where $V = \{v_1, \dots, v_N\}$ is a finite set of N elements [14]. In some cases it is possible that $u \neq v$, but $h(u) = h(v)$. Such situations are called collisions. The problem of collisions is solved by the technique called "chaining" [14]. The technique is to keep N distinct linked lists (or buckets) L_i one for each possible hash value. For $i \in \{1, \dots, N\}$ denote by B_i the set $\{x \in E^n / h(x) = v_i\}$. B_i 's are called blocks. The i -th list stores those vectors belonging to F , which have the same hash value, i.e. $L_i = \{x \in F / h(x) = v_i\}$ or in other way $L_i = B_i \cap F$. Hash coding scheme is called balanced, if $|B_i| = \frac{2^n}{N}$. The Elias algorithm [15] considers blocks B_i ordering them by their distances at vector x . Mention that we must have an efficient method to find all blocks $B_{j_1}, \dots, B_{j_{s(j)}}$ located at distance j from x if such blocks exist. After the step of ordering the algorithm examines the lists $L_{j_1}, \dots, L_{j_{s(j)}}$ one after the other by increase of j . Let the best match distance is denoted by δ . Due to $F \neq \emptyset$ initialisation of δ will happen on some step. Now, if the current values obey $\delta < j$ algorithm stops the work. All blocks with higher distances than δ at x do not need to be examined. In the reminder case $\delta \geq j$, examining nonempty list L_{j_k} algorithm can change the best match distance δ , also refreshing the current best match set, or the δ will remain unchanged and the current best match set will be updated. For balanced hash coding schemes it is proposed that the Elias algorithm may be optimal when the blocks B_i are isoperimetric sets [15, 16]. By the complexity of algorithm we mean the average number of examined lists over all files and queries, supposing that each vector $z \in E^n$ can independently appear in F with the same probability p . The pseudocode of the algorithm is brought below, where n is the word length, N is the number of blocks.

Elias algorithm

```

input  $x, F$ , comment:  $F \neq \emptyset$  integer  $\delta = \infty$ , comment: the current best match distance
set  $S = \emptyset$ , comment:  $S$  is the current set of vectors of  $F$  located at distance  $\delta$  from  $x$ 
integer  $j = -1$ ,
while( $j < \delta$ )
{  $j++$ ,
if( $s(j) \neq 0$ )
for(integer  $i = 0, i < s(j), i++$ )
{ if( $L_{j_i} \neq \emptyset$ ) comment: start examin the list  $L_{j_i}$ 
if( $\delta \leq d(x, L_{j_i})$ ) comment:  $\delta$  is unchanged

```

```

S = S ∪ (Oδn(x) ∩ Lji)
else { comment: δ is changed
S = Oδn(x) ∩ Lji,
δ = d(x, Lji) } } }
return S, comment: s = b(x, f), δ = d(x, f)

```

Now suppose we have an $[n, k, d]R$ code C . We define a hash function $h : E^n \rightarrow C$ associated to the code C in the following way:

$$h_C(x) = \{c_i / d(x, c_i) = d(x, C)\}. \quad (3)$$

As it follows from (3), $h_C(x)$ could be multivalued function because the blocks B_i are spheres of radius R and they can intersect. When the code C is perfect the mentioned blocks do not intersect and their union covers unit cube. The formula for complexity of algorithm is brought below for the case corresponding to Golay code. But as perfect codes exist in very simple cases [1, 2], we also consider hash functions associated to codes in some sense near to perfect codes. Such property have also the so called quasi-perfect codes [4–10]. Indeed the algorithm is proposed for balanced hash coding schemes, where different blocks B_i do not intersect, but we will also consider the algorithm for the case of intersecting blocks. In this case when blocks intersect we create the list in a similar way to the basic case and then these lists are also intersecting. Repeated element bring some redundancy (in terms of memory). The formal expression of complexity of algorithm is then brought for the particular case of two-error-correcting primitive BCH code of length $2^m - 1$ for odd m . To write a formula of complexity of the algorithm, for $x \in E^n$ let us consider the following Table:

x	p_1 F_1	p_2 F_2	...	$p_{2^{2^n}}$ $F_{2^{2^n}}$
B_1	α_{11}^x	α_{12}^x	...	$\alpha_{12^{2^n}}^x$
B_2	α_{21}^x	α_{22}^x	...	$\alpha_{22^{2^n}}^x$
...
B_{2^k}	$\alpha_{2^k 1}^x$	$\alpha_{2^k 2}^x$...	$\alpha_{2^k 2^{2^n}}^x$

In Table $F_1, \dots, F_{2^{2^n}}$ are all subsets of vertexes of unit cube and each F_i could be generated with the corresponding probability p_i . We will use the values α_{ij}^x putting them in the cells corresponding to block B_i and subset F_j , where

$$\alpha_{ij}^x = \begin{cases} 1, & \text{if } B_i \text{ considered in case of the set } F_j \text{ and vertex } x, \\ 0, & \text{otherwise.} \end{cases}$$

As we mention the complexity of algorithm will be represented as

$$\alpha(h_C) = \frac{1}{2^n} \sum_{x \in E^n} \sum_{1 \leq i \leq 2^k} \sum_{1 \leq j \leq 2^{2^n}} p_j \alpha_{ij}^x. \quad (4)$$

Let us denote $\Phi_x(B_i) = \sum_{1 \leq j \leq 2^{2^n}} p_j \alpha_{ij}^x$. As we can see $\Phi_x(B_i)$ is the probability that the block B_i will be considered by the algorithm when the vector x is requested. Then

$$\alpha(h_C) = \frac{1}{2^n} \sum_{x \in E^n} \sum_{1 \leq i \leq 2^k} \Phi_x(B_i).$$

It is easy to understand that for a fixed query x the block B_i will be examined, if the sphere $S_{d(x, B_i) - 1}^n(x)$ does not contain any vector belonging to F . In that case all blocks B_i such that

$d(x, B_i) \leq d(x, B_i) - 1$ will be examined. Let j vary over all possible distances between vector x and blocks B_i . Denote by $T_x(j)$ the number of blocks located at distance $\leq j$ from vector x then

$$\alpha(h_C) = \frac{1}{2^n} \sum_{x \in E^n} \sum_{0 \leq j \leq n} T_x(j) V(j), \quad (5)$$

where $V(j)$ denotes the probability that the nearest vector in F is located at distance j from x . Recall that [15] $V(j) = (1 - (1 - p)^{\binom{n}{j}})(1 - p)^{\sum_{i=0}^{j-1} \binom{n}{i}}$. As $d(x, c_i) = w(x + c_i)$ then the number of vectors located at distance i is equal to A_i^{x+C} . The sphere with centre c_i and radius R will be located at a distance $\leq j$ from vector x , if and only if $d(x, c_i) \leq j + R$. Therefore, $T_x(j) = \sum_{i=0}^{j+R} A_i^{x+C}$. Note that $A_i^{x+C} = 0$ when $i > n$.

As it is known the Golay code has four types of cosets [1] and each type can be obtained by some vector e_i of weight i , $i \in \{0, 1, 2, 3\}$. The number of cosets of minimum weight i is equal to $\binom{23}{i}$, and each coset contain 2^{12} vectors. Therefore, we get the following:

Proposition 1. For the Golay code the complexity of Elias algorithm is:

$$\begin{aligned} \alpha(h_{\Gamma_{23}}) = & \sum_{0 \leq j \leq 23} V(j) \sum_{i=0}^{j+3} \left(\frac{1}{2^{11}} A_i^{e_0+\Gamma_{23}} + \frac{23}{2^{11}} A_i^{e_1+\Gamma_{23}} + \right. \\ & \left. + \frac{253}{2^{11}} A_i^{e_2+\Gamma_{23}} + \frac{5819}{2^{11}} A_i^{e_3+\Gamma_{23}} \right). \end{aligned} \quad (6)$$

Proposition 1 gives the theoretical explanation of the experimental results, which are brought in [15]. As we mention for odd m two-error correcting BCH codes has four types of coset. Keeping in mind this and calculating the number of each type from (5) we get:

Proposition 2. For the two-error-correcting BCH code the complexity of Elias algorithm is:

$$\begin{aligned} \alpha(h_{BC_m}) = & \sum_{0 \leq j \leq 2^m - 1} V(j) \sum_{i=0}^{j+3} \left(\frac{1}{2^{2m}} A_i^{e_0+BC_m} + \frac{2^m - 1}{2^{2m}} A_i^{e_1+BC_m} + \right. \\ & \left. + \frac{(2^m - 1)(2^{m-1} - 1)}{2^{2m}} A_i^{e_2+BC_m} + \frac{2^{2m-1} + 2^{m-1} - 1}{2^{2m}} A_i^{e_3+BC_m} \right). \end{aligned} \quad (7)$$

Received 28.02.2013

REFERENCES

1. **Mac-Williams F.J., Sloane N.J.** The Theory of Error-Correcting Codes. Amsterdam: N.-H. Mathematical Library, 1977, 762 p.
2. **Zinoviev V.A., Leont'ev V.K.** The Nonexistence of Perfect Codes Over Galois Fields. // Problems of Control and Info. Theory, 1973, v. 2, № 2, p. 123–132.
3. **Bassalygo L.A., Zaitsev G.V., Zinoviev V.A.** Uniformly Packed Codes. // Problemi Peredachi Informatsii, 1974, v. 10, № 1, p. 9–14 (in Russian).
4. **Cohen G., Honkala I., Litsyn S., Lobstein A.** Covering codes. Amsterdam: N.-H. Mathematical Library, 1997, 542 p.
5. **Baicheva T., Bouyukliev I., Dodunekov S.** Binary and Ternary Linear Quasi-Perfect Codes with Small Dimensions. // Transactions of Information Theory, 2008, v. 54, № 9, p. 4335–4339.

6. **Gabidulin E.M., Davydov A.A., Tombak L.M.** Linear Codes with Covering Radius 2 and Other New Covering Codes. // Transactions of Information Theory, 1991, v. 37, № 1, p. 219–224.
7. **Davydov A.A., Tombak L.M.** Quasiperfect Linear Binary codes with Distance 4 and Complete Caps in Projective Geometry. // Problemi Peredachi Informatsii, 1989, v. 25, № 4, p. 265–275 (in Russian).
8. **Davydov A.A., Drozhzhina-Labinskaya A.Yu.** Constructions, Families, and Tables of Binary Linear Covering Codes. // Transactions of Information Theory, 1994, v. 40, № 4, p. 1270–1279.
9. **Etzion T., Mounits B.** Quasi-Perfect Codes with Small Distance. // Transactions of Information Theory, 2005, v. 51, № 11, p. 3938–3946.
10. **Etzion T., Greenberg G.** Constructions for Perfect Mixed Codes and Other Covering Codes. // Transactions of Information Theory, 1993, v. 39, № 1, p. 209–214.
11. **Gorenstein D., Peterson W., Zierler N.** Two Error-Correcting Bose-Chaudhuri Codes are Quasi Perfect. // Information and Control, 1960, v. 3, № 3, p. 291–294.
12. **Kasami T., Lin S., Peterson W.** Some Results on the Weight Distributions of BCH Codes // Transactions of Information Theory, 1966, v. 12, № 2, p. 274–277.
13. **Charpin P.** Weight Distributions of Cosets of Two-Error-Correcting BCH Codes, Extended or Not. // Transactions of Information Theory, 1994, v. 40, № 5, p. 1425–1442.
14. **Knuth D.E.** The Art of Computer Programming. V. 3. Sorting and Searching. Massachusetts: Addison-Wesley, 1998, 780 p.
15. **Rivest R.L.** On the Optimality of Elias's Algorithm for Performing Best-Match Searches. // Information Processing, 1974, p. 678–681.
16. **Aslanyan L.H.** The Discrete Isoperimetry Problem and Related Extremal Problems for Discrete Spaces. // Problemi Kibernetiki, 1979, v. 36, p. 85–128 (in Russian).