

On the Optimality of a Hash-Coding Type Search Algorithm

L. H. Aslanyan

Institute for Informatics and Automation
Problems of NAS RA

Yerevan, Armenia

e-mail: lasl@sci.am

H. E. Danoyan

Institute for Informatics and
Automation Problems of NAS RA

Yerevan, Armenia

e-mail: haykdanoyan@yandex.ru

ABSTRACT

The algorithmic problem of nearest neighbors search and the early proposed hash-coding algorithms (Elias algorithm, tree based algorithms) are considered. The problem of optimality characterization of the algorithm - depending on the geometrical structure of blocks in a hash-coding scheme is considered. As a supportive technique for this the structure of consecutive neighbourhood layers to the subsets of binary cube is considered. It is proved that the function representing the cardinalities of k^{th} neighborhoods of a set depending on k is concave. Using this fact and the known properties of isoperimetric sets the difference between the mentioned functions (for fixed argument) corresponding to isoperimetric set and an arbitrary set of the same cardinality can be estimated. It is proved that the algorithm is optimal when the mentioned blocks are isoperimetric sets.

Keywords

Best match, nearest neighbors, hash-coding, Elias algorithm, k^{th} neighborhood, standard placement.

1. INTRODUCTION

Let $E = \{0,1\}$. Consider Cartesian degree E^n which is known as the set of vertexes of n -dimensional unit cube. For $x, y \in E^n$ denote by $d(x, y)$ the Hamming distance between vectors x and y . Denote by $w(x)$ the weight of vector x , i.e. $w(x) = d(0^n, x)$. 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 neighbourhood / shell of radius r i.e. $O_r^n(x) = \{y \in E^n / d(x, y) = r\}$. For an $x \in E^n$ and $A \subseteq E^n$ we define the distance $d(x, A)$ as $d(x, A) = \min_{a \in A} d(x, a)$. For $A \subseteq E^n$ denote by $A^{(k)}$ the k^{th} neighborhood of A , i.e. $A^{(k)} = \{x / d(x, A) = k\}$.

The algorithmic problem of finding the set of all nearest neighbors to a given vector from a given set is known [1]. More precisely, let us we have a nonempty subset $F \subseteq E^n$ and a vector x (query element or query vector). Denote by F_x the set of nearest neighbors to x from F i.e.

$$F_x = \{y \in F / d(x, y) = d(x, F)\}.$$

The problem is for a given F and x to find in an optimal way the set F_x . We can find the mentioned set by calculating for each $y \in F$ the corresponding distances to x but in many cases, for large F it will make the problem practically unsolvable. So, it is required from algorithm to consider as

few as possible elements from F in constructing the nearest neighbours. One algorithm solving the mentioned problem is known as Elias algorithm [1]. The main goal of this paper is to provide the proof of the conjecture about optimality of the mentioned algorithm concerning the geometrical shape of code blocks[1].

The paper is organized as follows: In section 2 the descriptions are brought related to the hash-coding schemas and algorithms. In section 3 some auxiliary results and definitions are brought about cardinalities of neighborhoods of a cube subset. Finally, in section 4 the theorem is proved about the optimality of the algorithm.

2. ELIAS ALGORITHM

At first, let us describe a particular schema of representation of subset F , called hash-coding schema. Hash function is defined as a function $h: E^n \rightarrow V$, where V is a finite set of N elements, $V = \{v_1, \dots, v_N\}$ [2]. We consider the case, that $V = E^m$, and $m \leq n$. In connection to hash function, F is represented as a union of N distinct linked lists (or buckets) L_i , $i = 1, \dots, N$ 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 words $L_i = F \cap B_i$. Hash coding scheme is called balanced if $|B_i| = 2^n / N$, for $i = 1, \dots, N$.

The Elias algorithm [1] considers blocks B_i of the model, 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 these ordering the algorithm examines the lists $L_{j_1}, \dots, L_{j_s(j)}$ one after the other by increase of j_t . Let the best match distance (also its current value in algorithm) be denoted by δ . Due to $F \neq \emptyset$ initialization of δ will happen on some step and the final algorithmic value will present the real minimal distance δ . 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_i} 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.

By the complexity of algorithm we mean the average number of examined lists over all files and queries, supposing that

- I. Each vector $x \in E^n$ equally likely can be requested.
- II. Each vector $z \in E^n$ independently appears in F with the same probability p . This gives a specific probabilistic distribution over the set of subsets of E^n appearing as the target set F .

The pseudocode of the algorithm is brought below:

Elias Algorithm /* n is the word length, N is the number of blocks */
input x, F , /* $F \neq \emptyset$ */
integer $\delta = \infty$, /* the current best match distance */
set $S = \emptyset$, /* S is the current set of vectors of F located at distance δ from x */
integer $j = -1$, /* current distance of blocks under consideration from x */
while ($j < \delta$)
{
 $j++$,
if ($s(j) \neq 0$) /* $s(j)$ is the number of blocks in distance j from x */
for (integer $i=0, i < s(j), i++$)
{
If ($L_{ji} \neq \emptyset$) /* start examine the list L_{ji} , i -th list with j distance block */
If ($\delta \leq d(x, L_{ji})$) /* δ is unchanged */
 $S = S \cup (O_\delta^n(x) \cap L_{ji})$ /* $O_\delta^n(x)$ is the δ neighborhood of x */
else
{
 $S = O_\delta^n(x) \cap L_{ji}$, /* δ is changed */
 $\delta = d(x, L_{ji})$
}
}
}
return S , comment: $S = F_x, \delta = d(x, F)$

For balanced hash coding schemes it is proposed that the algorithm may be optimal when the blocks B_i are isoperimetric sets [1],[3]. Our main goal is to prove the mentioned conjecture.

3. ON FUNCTION OF CARDINALITIES OF NEIGHBORHOODS OF A SET

For $A \subseteq E^n$ let us denote

$$\Gamma_A(k) = |A^{(k)}|, k \in \{0, \dots, n\}.$$

It is obvious, that $\sum_{k=0}^n \Gamma_A(k) = 2^n$. For a one-point set, $A = \{x\}$, values $\Gamma_A(k)$ are given by combinations $\binom{n}{k}$. For larger sets starting from some point values $\Gamma_A(k)$ become equal to 0. For a one-point set it is evident that $\Gamma_A(k)$ is simply increasing and then decreasing at the point $n/2$. The similar behaviour of the function $\Gamma_A(k)$ is not studied for arbitrary $A \subseteq E^n$ and our nearest goal is this investigation.

Lemma 3.1. For arbitrary subset A the function $\Gamma_A(k)$ has one or two maximum points on the domain $\{1, \dots, n\}$ and in case when the maximum points are two then they differ by 1.

Proof. Let us prove the lemma by induction. For $n = 3$ one can check that the lemma takes place. Now let the lemma be true in case $\leq n - 1$ and prove it for the case of n . Let us fix the i^{th} coordinate and consider the partition of unit cube by coordinate x_i . Let us denote the corresponding subcubes by $E_{x_i=0}^n$ and $E_{x_i=1}^n$. In connection to this, denote the subsets of the set A belonging to $E_{x_i=0}^n$ and $E_{x_i=1}^n$ by $A_{x_i=0}$ and $A_{x_i=1}$ correspondingly. Let us denote $P_0 = A_{x_i=0} \cup A_{x_i=0}^{(1)} \cup A_{x_i=1}$ and $P_1 = A_{x_i=1} \cup A_{x_i=1}^{(1)} \cup A_{x_i=0}$. Note that P_0 and P_1 are subsets of $n - 1$ dimensional unit cube. It can be shown that (Figure 1)

$$\Gamma_A(k) = \Gamma_{P_0}(k-1) + \Gamma_{P_1}(k-1) \quad (3.1)$$

but for the sake of space here we omit the detail proof of the claim.

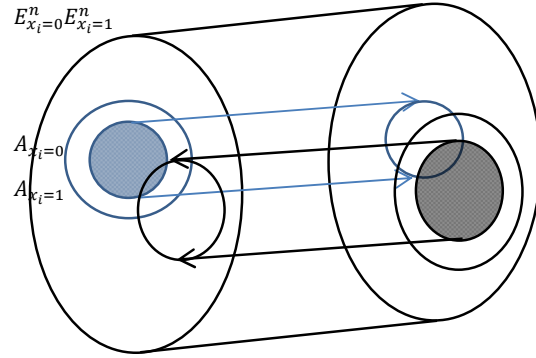


Figure 1.

As P_0 and P_1 are subsets of $n - 1$ dimensional unit cube, then for functions $\Gamma_{P_0}(k)$ and $\Gamma_{P_1}(k)$ the lemma is true by the induction hypothesis. Now let k_0 be the smallest integer such that one of the mentioned functions (without losing generality assume that it is the $\Gamma_{P_0}(k)$) decreases. Then from the definition of the set P_1 it follows that at most at point $k_0 + 1$ the function $\Gamma_{P_1}(k)$ also decreases. The lemma follows purporting this and equation (3.1).

Now let us have a number $0 \leq a \leq 2^n$ and a point $c \in E^n$. We define a set of cardinality a called standard placement [3]. At first we define a linear order on E^n by the following way [1],[3]. Let $u, v \in E^n$, then:

- I. if $d(u, c) < w(t(v, c))$ then $u < v$,
- II. if $d(u, c) = d(v, c)$ then $u < v$ if there exists $i, 1 \leq i \leq n$ such that $u_j = v_j, j = 0, \dots, i - 1$ and $u_i + c_i = 1, v_i + c_i = 0$.

The standard placement with centre c and cardinality a is denoted by $L_a^n(c)$. $L_a^n(c)$ is the set of first a elements by means of linear order defined above [1],[3].

Theorem 3.1 [1],[3]. For $A \subseteq E^n, |A| = a$ and $c \in E^n$ the following takes place:

$$\Gamma_{L_a^n(c)}(1) \leq \Gamma_A(1). \quad (3.2)$$

This result is the origin of our investigation. It is known as the base Discrete Isoperimetry Theorem. More properties of Discrete Isoperimetry were investigated, see [3],[4], but this research concerns only the set and its 1-neighbourhood. Consecutive neighbourhoods raise more questions, and here our target is understanding the functional behaviour of the series of neighbourhoods, and

their sizes in particular. Now let us have an arbitrary subset A , $|A| = a$ and the standard placement $L^n_a(c)$.

Lemma 3.2. There exists an integer s ($1 \leq s \leq n$) such that

$$\Gamma_{L^n_a(c)}(k) \leq \Gamma_A(k), \text{ when } k \leq s \quad (3.3)$$

and

$$\Gamma_{L^n_a(c)}(k) \geq \Gamma_A(k), \text{ when } k > s. \quad (3.4)$$

The proof technique is similar to the proof of Lemma 3.1 and is omitted due to the space limitation. For geometrical interpretation of lemma 3.2 see figure 2.

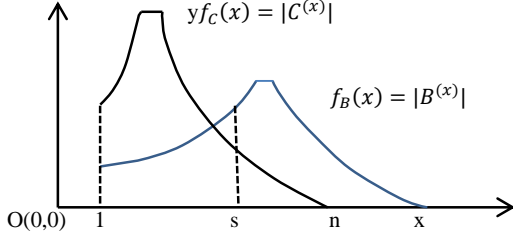


Figure 2.

4. OPTIMALITY OF THE HASH-CODING TYPE SEARCH ALGORITHM

Let us consider the following problem: which property must satisfy blocks of balanced hash-coding schema, so the algorithm to be optimal.

Theorem 4.1. The algorithm is optimal, when the blocks B_i , $i = 1, \dots, N$, of hash-coding schema are standard placements (spheres in particular cases).

Proof. As it is known [1], the complexity of the algorithm is

$$\alpha(h) = \sum_{i=1}^N \Phi(B_i), \quad (4.1)$$

where by $\Phi(B_i)$ is denoted the average probability aver all query vectors that the block B_i will be considered. It is known [1], that

$$\Phi(B_i) = \frac{1}{2^n} \sum_{j=0}^n |B_i^{(j)}| \Psi(n, p, j-1), \quad (4.2)$$

where $\Psi(n, p, j-1)$ is the probability that the sphere of radius $j-1$ does not contain points belonging to subset F . Now let B be a standard placement and C be an arbitrary set, such that $|B| = |C|$. Let us consider the following difference

$$\Phi(C) - \Phi(B) = \quad (4.3)$$

$$= \frac{1}{2^n} \sum_{i=1}^n (|C^{(i)}| - |B^{(i)}|) \cdot \Psi(n, p, i-1)$$

From theorem 3.2 and lemma 3.2 it follows that there exists an integers, such that

$$\Phi(C) - \Phi(B) =$$

$$\begin{aligned} &= \frac{1}{2^n} \sum_{i=1}^s (|C^{(i)}| - |B^{(i)}|) \cdot \Psi(n, p, i-1) \\ &+ \frac{1}{2^n} \sum_{i=s+1}^n (|C^{(i)}| - |B^{(i)}|) \cdot \Psi(n, p, i-1) \\ &\geq \frac{1}{2^n} \sum_{i=m}^s (|C^{(i)}| - |B^{(i)}|) \cdot \Psi(n, p, s-1) - \\ &- \frac{1}{2^n} \sum_{i=s+1}^n (|B^{(i)}| - |C^{(i)}|) \cdot \Psi(n, p, s) \end{aligned} \quad (4.3)$$

As $\Psi(n, p, s) \leq \Psi(n, p, s-1)$ and $\sum_{i=1}^s (|C^{(i)}| - |B^{(i)}|) = \sum_{i=s+1}^n (|B^{(i)}| - |C^{(i)}|)$, then from (4.3) follows that

$$\Phi(C) - \Phi(B) \geq 0.$$

CONCLUSION

The nearest neighbors search is an important part of many procedures such as in Pattern Recognition, Natural Language Error Correction, and others. Two types of algorithms are well known in this area: the tree based algorithms and the dynamic programming style algorithms. Our interest is devoted to the second one. First of all we use the hypotheses that the optimality of this type algorithms is related to the well known Discrete Isoperimetry problem. Then using the known results in this area and giving some generalization we obtain combinatorial postulates that help to prove the optimality of these algorithms. Results achieved are interesting for both – the Discrete Isoperimetry domain and in Divide and Conquer type search. The complementary domain of investigation that considers the issues how to partition the discrete space to the Isoperimetry type blocks is also under consideration but this is a separate broad research direction that is addressed [5], [6].

REFERENCES

- [1] R. L. Rivest, "On the optimality of Elias's algorithm for performing best-match searches", Information Processing, pp. 678-681, 1974.
- [2] D. E. Knuth, "The Art of Computer Programming, vol. 3 / Sorting and Searching", Eddison-Wesley, p. 800, 1998.
- [3] L. H. Aslanyan, "The Discrete Isoperimetry Problem and Related Extremal Problems for Discrete Spaces", Problemi Kibernetiki, pp. 85-128, 1979.
- [4] L. H. Harper, Global Methods for Combinatorial Isoperimetric Problems, Cambridge University Press, 2010, 250p.
- [5] L. H. Aslanyan, H. E. Danoyan, Complexity of Elias algorithm based on codes with covering radius three, Proceedings of the Yerevan state university, 2013 N°1, pp. 44-50.
- [6] L. H. Aslanyan, H. E. Danoyan, Complexity of Elias algorithm based on Hamming and extended Hamming codes, Reports of NAS RA, vol. 113, no. 2, pp. 151-158, 2013.