Astérisque

GILLES COHEN GÉRARD ZÉMOR Subset sums and coding theory

Astérisque, tome 258 (1999), p. 327-339 <http://www.numdam.org/item?id=AST 1999 258 327 0>

© Société mathématique de France, 1999, tous droits réservés.

L'accès aux archives de la collection « Astérisque » (http://smf4.emath.fr/ Publications/Asterisque/) implique l'accord avec les conditions générales d'utilisation (http://www.numdam.org/conditions). Toute utilisation commerciale ou impression systématique est constitutive d'une infraction pénale. Toute copie ou impression de ce fichier doit contenir la présente mention de copyright.

\mathcal{N} umdam

Article numérisé dans le cadre du programme Numérisation de documents anciens mathématiques http://www.numdam.org/ Astérisque 258, 1999, p. 327–339

SUBSET SUMS AND CODING THEORY

by

Gérard Cohen & Gilles Zémor

Abstract. — We study some additive problems in the group $(\mathbb{Z}_{2\mathbb{Z}})^r$. Our purpose is to show how those problems are closely related to coding theory. We present some relevant classical coding techniques and make use of them to obtain some original contributions.

1. Introduction

Let G denote the group \mathbf{F}^r where $\mathbf{F} = \{0, 1\}$ stands for the additive group with two elements. Let S be a generating set of G. For any positive integer *i*, denote by S^i the set of sums of *i* distinct elements of S. Set $S^0 = \{0\}$ and for any set I of non-negative integers, let $S^I = \bigcup_{i \in I} S^i$. Let us denote by $\rho(S)$ the smallest integer *t* such that any element of G can be expressed as a sum of *t* or less elements of S, i.e. such that

$$G = S^{[0,t]}.$$

Let us denote by d(S) the smallest integer *i* such that 0 can be expressed as a sum of *i* distinct elements of *S*, i.e. let d(S) - 1 be the largest *t* such that

 $0 \not\in S^{[1,t]}.$

We wish to focus on the following 'additive' problems.

Problem 1. — For given r and t, find the smallest s such that $|S| \ge s$ implies $\rho(S) \le t$.

Problem 2. — For given r and t, find the largest s such that $|S| \leq s$ implies $\rho(S) \geq t$.

Problem 3. — For given r and d, find the smallest s such that $|S| \ge s$ implies $d(S) \le d$.

¹⁹⁹¹ Mathematics Subject Classification. — 94B05, 94B75, 05C25, 05C50, 11P99. Key words and phrases. — Coding theory, additive theory.

Those three problems can be expressed as problems in coding theory. Indeed, problems 2 and 3 are classical coding problems of which we shall give a short self-contained presentation for the non specialist. Problem 1, although less known to coding theorists, is also amenable to coding techniques, and we shall present original contributions to it and also to the following generalisation of problem 3.

Problem 4. — Given r and an arbitrary set of integers I, find the smallest s such that $|S| \ge s$ implies $0 \in S^I$.

2. Coding-theoretic formulation of problems 1-4

What coding theorists call a (binary) *linear code* of length n is simply a subspace of the vector space \mathbf{F}^n . Let S be a generating set of \mathbf{F}^r with |S| = n. There is an important linear code C(S) associated to S whose coding-theoretic properties reflect the additive properties of S. To obtain it let s_1, \ldots, s_n be any ordering of its elements that we shall write as column vectors. Consider the $r \times n$ matrix $\mathbf{H} = [s_1 \ldots s_n]$ and the associated function

$$\sigma: \mathbf{F}^n \to G = \mathbf{F}^r$$
$$\mathbf{x} = (x_1 \dots x_n) \mapsto \sigma(\mathbf{x}) = \mathbf{H}^t \mathbf{x}$$

Define C(S) to be the set of vectors \mathbf{x} of \mathbf{F}^n such that $\sigma(\mathbf{x}) = 0$. When defining such a code C(S) associated to a set S we shall usually not specify which ordering s_1, \ldots, s_n we are choosing because the properties of C(S) that interest us are independent of it. To help distinguish between the two structures $G = \mathbf{F}^r$ and \mathbf{F}^n , we shall use plain letters to denote elements of G and bold letters to denote vectors of \mathbf{F}^n : furthermore, since the vector space structure of \mathbf{F}^n will be used rather more heavily than that of G, we shall systematically refer to elements of \mathbf{F}^n as vectors. C(S) (or simply C when there is little ambiguity) is a subspace of \mathbf{F}^n of dimension k = n - r. Its elements are referred to as codewords. H is called a parity-check matrix of C, and for any vector $\mathbf{x} \in \mathbf{F}^n$, $\sigma(\mathbf{x})$ is called the syndrome of \mathbf{x} . Two vectors $\mathbf{x} = (x_1 \dots x_n)$ and $\mathbf{y} = (y_1 \dots y_n)$ of \mathbf{F}^n are said to be orthogonal if

$$\sum_{i=1}^{n} x_i y_i = 0$$

where computations are performed in **F**. If C is a linear code of \mathbf{F}^n of dimension k, then the set C^{\perp} of vectors orthogonal to C is a linear code of dimension n - k. Any matrix **H** whose rows are independent vectors orthogonal to C make up a parity-check matrix of C.

Remark. — Not every code C need be a code C(S) for some set S. This is because not every code has a parity-check matrix with distinct columns.

Coding theorists regard \mathbf{F}^n as a metric space, i.e. endowed with the Hamming distance $d(\cdot, \cdot)$:

$$\begin{array}{rccc} \mathbf{F}^n \times \mathbf{F}^n & \to & [0,n] \\ (\mathbf{x},\mathbf{y}) & \mapsto & d(\mathbf{x},\mathbf{y}) \end{array}$$

where $d(\mathbf{x}, \mathbf{y})$ is defined as the number of coordinates where \mathbf{x} and \mathbf{y} differ. The *minimum distance* d(C) of a code C is the smallest distance between a pair of distinct codewords,

$$d(C) = \min_{\substack{\mathbf{x}, \mathbf{y} \in C \\ \mathbf{x} \neq \mathbf{y}}} d(\mathbf{x}, \mathbf{y}).$$

Note that d(C) is also the minimum distance $d(\mathbf{x}, \mathbf{0})$ between the **0** vector and any non-zero codeword \mathbf{x} : this is because $d(\cdot, \cdot)$ is invariant by translation and C is an additive subgroup. The integer $d(\mathbf{x}, \mathbf{0})$ is called the *weight* of \mathbf{x} and denoted by $w(\mathbf{x})$. The classical parameters of a linear code C are usually denoted by [n, k, d] and refer respectively to its length, dimension and minimum distance.

Another classical parameter of a code C is its covering radius $\rho(C)$: it is the maximum distance between a vector of \mathbf{F}^n and the code C, i.e.

$$\rho(C) = \max_{\mathbf{x} \in \mathbf{F}^n} d(\mathbf{x}, C)$$

where $d(\mathbf{x}, C) = \min_{\mathbf{c} \in C} d(\mathbf{x}, \mathbf{c})$.

Given a vector $\mathbf{x} = (x_1 \dots x_n)$ of \mathbf{F}^n , it is common to define its *support* by $supp(\mathbf{x}) = \{i, x_i = 1\}$. The syndrome of \mathbf{x} can therefore be written as

$$\sigma(\mathbf{x}) = \sum_{i \in supp(\mathbf{x})} s_i$$

where the sum is computed in \mathbf{F}^r . It is now clear that the minimum distance of C equals the minimum cardinality of a subset I of S such that $\sum_{i \in I} s_i = 0$. In particular we have :

Remark. — For any code C, there exists a set S not containing 0 such that C = C(S) if and only if $d(C) \ge 3$.

Similarly, it is readily checked that the covering radius of C is the smallest number of additions necessary to generate every non-zero element of \mathbf{F}^r with elements of S. Summarizing,

Proposition 2.1. — The correspondence $S \to C(S)$ is such that

$$d(S) = d(C(S))$$

$$\rho(S) = \rho(C(S)).$$

The above correspondence transforms problems of an additive nature into *packing* and *covering* problems in a metric space. In particular, we see that problem 3 is equivalent to the fundamental problem of coding theory, namely determine the largest possible minimum distance of a linear code of length n and dimension k. There are several classical bounds relating n, k and d. Let us mention two simple bounds that we shall make use of later on.

Proposition 2.2 (Hamming bound). — Any [n, n - r, d] code satisfies

$$\sum_{i=0}^{\lfloor (d-1)/2 \rfloor} \binom{n}{i} \le 2^r.$$

Proof. — Since any vector $\mathbf{x} \in \mathbf{F}^n$ of weight $\leq d-1$ satisfies $\sigma(\mathbf{x}) \neq 0$, then all vectors with weight at most $\lfloor (d-1)/2 \rfloor$ must have distinct syndromes.

Using classical estimates for binomial coefficients, the Hamming bound states, asymptotically, that any $[n, nR, n\delta]$ code satisfies

(1)
$$R \le 1 - h(\delta/2) + o(1)$$

where $h(x) = -x \log_2 x - (1-x) \log_2(1-x)$ denotes the binary entropy function.

Proposition 2.3 (Varshamov-Gilbert bound). — Let n and r be given. There exists an [n, n-r, d] code whenever

$$\sum_{i=0}^{d-1} \binom{n-1}{i} < 2^r.$$

Proof. — We construct inductively a parity-check matrix of such a code. Suppose constructed an $r \times i$ matrix \mathbf{H}_i such that any d-1 columns are linearly independent. They are at most N_i distinct linear combinations of columns involving at most d-2 terms, with

$$N_i = \sum_{j=1}^{d-2} \binom{n}{j}.$$

If $N_i < 2^r - 1$, then a nonzero element of $G = \mathbf{F}^r$ can be added to the set of columns of \mathbf{H}_i to yield an $r \times (i + 1)$ matrix \mathbf{H}_{i+1} with the property that any d - 1 of its columns are linearly independent; equivalently \mathbf{H}_{i+1} is the parity-check matrix of a code of minimal distance $\geq d$.

Asymptotically, the Varshamov-Gilbert bound reads: there exist $[n, nR, n\delta]$ codes with

(2)
$$R \ge 1 - h(\delta) + o(1).$$

There is no known better asymptotic lower bound on R = k/n. Let us just mention the most powerful upper bound on R due to McEliece, Rodemich, Rumsey, and Welch (see e.g. [10]) for a proof):

Proposition 2.4. — Any $[n, nR, n\delta]$ code satisfies

(3)
$$R \le h\left(\frac{1}{2} - \sqrt{\delta(1-\delta)}\right) + o(1).$$

Note that the Varshamov-Gilbert bound is not really constructive (the complexity of constructing a parity-check matrix for such codes is exponential in the length n). There are no known constructions of codes achieving the Varshamov-Gilbert bound for growing n and fixed R, 0 < R < 1. There are, however, good constructions of codes with fixed d and growing n. We give a very short presentation of such codes, to which we shall refer later on.