

Linear Boolean classification, coding and “the critical problem”

Emmanuel Abbe
Princeton University
eabbe@princeton.edu

Noga Alon
Tel Aviv University and
Institute for Advanced Study, Princeton
nogaa@tau.ac.il

Afonso S. Bandeira
Princeton University
ajs@math.princeton.edu

Abstract—This paper considers the problem of linear Boolean classification, where the goal is to determine in which set, among two given sets of Boolean vectors, an unknown vector belongs to by making linear queries. Finding the least number of queries is formulated as determining the minimal rank of a matrix over $GF(2)$ whose kernel does not intersect a given set S . In the case where S is a Hamming ball, this reduces to finding linear codes of largest dimension. For a general set S , this is an instance of “the critical problem” posed by Crapo and Rota in 1970, open in general. This work focuses on the case where S is an annulus. As opposed to balls, it is shown that an optimal kernel is composed not only of dense but also of sparse vectors, and the optimal mixture is identified in various cases. These findings corroborate a proposed conjecture that for an annulus of inner and outer radius nq and np respectively, the optimal relative rank is given by the normalized entropy $(1-q)H(p/(1-q))$, an extension of the Gilbert-Varshamov bound.

I. INTRODUCTION

A. Linear coding

One of the fundamental problems of coding theory is to identify the largest dimension of a binary code with a given length and distance. This means to identify the largest cardinality of a subset of \mathbb{F}_2^n whose elements are at distance at least d from each other. This is a well-known open problem in general. Even for the case of a linear code, i.e., a subspace of \mathbb{F}_2^n , the problem is open. From the parity-check matrix viewpoint, constructing a linear code of distance d is equivalent to constructing a matrix M such that Mx allows to recover x for all x having weight at most $\lfloor (d-1)/2 \rfloor$, equivalently, to construct a matrix M of least rank such that

$$Mx \neq Mx', \quad \forall x, x' \in B(0^n, s), x \neq x', \quad (1)$$

where

$$B(0^n, s) = \{x \in \mathbb{F}_2^n : w(x) \leq s\}, \quad (2)$$

$$s = \lfloor (d-1)/2 \rfloor, \quad (3)$$

is the Hamming ball of centre 0^n and radius s , and where $w(x)$ denotes the Hamming weight of x . Note that for a fixed s and n , the least rank of matrices satisfying property (1) is a finite integer, we denote it by $m^*(2s, n)$. The factor 2 will be justified later. Finding m^* for general values is a difficult problem as mentioned previously, and even in the asymptotic regime of s, n diverging with a fixed ratio $s/n = \delta/2$, the problem is still open. In fact, it is believed by some (e.g., Goppa’s conjecture) that the answer is given by

$$m^*(\delta n, n) = nH(\delta) + o(n), \quad (4)$$

the Gilbert-Varshamov (GV) bound [5], [9], where

$$H(\delta) = -\delta \log_2 \delta - (1-\delta) \log_2 (1-\delta), \quad (5)$$

if $\delta \in [0, 1/2]$ and $H(\delta) = 1$ if $\delta \in (1/2, 1]$.

It is not difficult to establish this bound, with a greedy algorithm or with a probabilistic argument, but it has not been improved since the 50’s in the asymptotic regime, nor has it been proved to be tight.

Linear codes are particularly interesting for several reasons. Their encoding complexity is reduced from exponential (in the worst case) to quadratic in the blocklength (specifying a basis). Moreover, most of the codes studied in the literature and used in applications with efficient decoding algorithms are linear. There are also other interesting features of linear codes, such as their duality with linear source codes. The parity-check matrix of a linear code can be viewed as a linear source compressor, for a source distribution (or a source model in the worst-case setting) equivalent to the error distribution (or model). In particular, if the source model is given by the k -sparse sequences, i.e., binary sequences with at most k ones, then the optimal compression dimension, assuming the Gilbert-Varshamov bound to be tight, is given by $nH(2k/n) + o(n)$, where the first term is approximately $2k \log n/k$, when k/n is small.

B. Linear Boolean classification

Given two disjoint classes $S_1, S_2 \subseteq \{0, 1\}^n$, we are interested in constructing a linear map $M : \mathbb{F}_2^n \rightarrow \mathbb{F}_2^m$ such that for $x \in S_1 \cup S_2$, Mx allows to determine if x belongs to S_1 or S_2 , i.e., $MS_1 \cap MS_2 = \emptyset$. We refer to this problem as linear Boolean classification. Our goal is to identify the least dimension m for which linear classification is possible when the two classes are given.

Consider the following toy example. The unknown vector $x \in \mathbb{F}_2^n$ represents the pixels of a black and white image. Assume that there are two models for the image, the dark and bright images, modelled by the two sets $S_1 = B(1^n, s)$ and $S_2 = B(0^n, s)$, where $s < n/2$. This means that the bright images are all vectors with at most s black pixels, whereas the dark images are all vectors with at most s white pixels. Our goal is to identify the least dimension of a linear projection for which we can still classify if an image is dark or bright. The reader can easily check that in this case the least dimension is $2s + 1$, as further discussed in the next section. This shows that linear classification, as compared to linear compression, can be simpler to address for Hamming balls. The reason behind this difference will be clarified in the next section.

Note that we do not specify how to decide among the two classes by accessing $y = Mx$. A general approach is to find a solution x_0 of $y = Mx$, and to then verify if $x_0 + \ker(M)$ intersects S_1 or S_2 . This may be of course computationally costly, but our goal in this paper is only to identify the least rank of M for which the intersection always happens only with one of the two sets. The computational efficiency would be an interesting problem to pursue.

II. PRELIMINARY RESULTS AND “THE CRITICAL PROBLEM”

A. Linear coding for general models

As expressed in (1), linear coding can be viewed as constructing flat matrices which are injective for sequences constrained to have a bounded number of ones. This concerns linear coding for the traditional Hamming ball model. One can consider more general models, in which case the injectivity property (1) needs to be guaranteed for vectors x belonging to a specified set $S \subseteq \mathbb{F}_2^n$. From now on, we define linear codes by means of parity-check matrices.

Definition 1. A linear code $M : \mathbb{F}_2^n \rightarrow \mathbb{F}_2^m$ can compress losslessly a source model $S \subseteq \mathbb{F}_2^n$ (or can correct the error patterns in $S \subseteq \mathbb{F}_2^n$), if

$$Mx \neq Mx', \quad \forall x, x' \in S, x \neq x', \quad (6)$$

i.e., if M is S -injective.

Definition 2. For a given set $S \subseteq \mathbb{F}_2^n$, we define the linear compression dimension of S by

$$m^*(S) = \min_{M \in S\text{-injective}} \text{rank}(M), \quad (7)$$

and the linear compression rate of S by $m^*(S)/n$.

Note that $m^*(S)$ can be equivalently defined by

$$m^*(S) = n - \max_{V: (V \setminus \{0\}) \cap (S+S) = \emptyset} \dim(V), \quad (8)$$

where V denotes a subspace of \mathbb{F}_2^n . In other words, we need to find the largest dimension of a subspace avoiding $(S+S) \setminus \{0\}$. Throughout the paper, the sum of two sets is defined by $S_1 + S_2 = \{s_1 + s_2 : s_1 \in S_1, s_2 \in S_2\}$. The problem of finding the largest dimension of a vector space which does not intersect a given subset of \mathbb{F}_q^n (where q is a power of a prime) is known as “the critical problem” and was posed¹ by Crapo and Rota in [3]. Even for $q = 2$, it is an open problem for arbitrary sets.

Note also that if one is allowed to use a non-linear map M , required to be S -injective, the “dimension” of M can be as low as $\log_2(|S|)$, but not lower. Quotes are used on “dimension” since the map is non linear and since a priori $\log_2(|S|)$ may not be an integer. For linear maps, we then have

$$\log_2(|S|) \leq m^*(S). \quad (9)$$

One can also obtain the following upper-bound with a probabilistic argument.

Lemma 1. For any $S \subseteq \mathbb{F}_2^n$,

$$\log_2(|S|) \leq m^*(S) \leq \lfloor \log_2(|S+S| - 1) \rfloor + 1. \quad (10)$$

A proof of this simple lemma is available in [1]. For Hamming balls, the above bounds are equivalent to the Hamming and Gilbert-Varshamov bounds. Note that if S is a subspace, then the bounds match and are equal to $\log_2(|S|)$.

B. Linear coding for subspaces

The lower bound (9) is clearly achieved if S is a subspace of \mathbb{F}_2^n , i.e., if $S+S = S$, using for M the projection on S . One may ask for what kind of set S are the two bounds in the lemma matching in the asymptotic regime, i.e., up to $o(n)$. This is equivalent to asking for the doubling constant of S to be sub-exponential, i.e.,

$$\frac{|S+S|}{|S|} = 2^{o(n)}. \quad (11)$$

¹In fact, an even more general formulation is proposed in [3]

One may expect that this holds only if the set S is closed to a subspace in some sense. In fact, this is related to the Polynomial Freiman-Ruzsa conjecture (see [6]):

Conjecture 1. [Polynomial Freiman-Ruzsa conjecture] *If S has a doubling constant at most K , then S is contained in the union of $K^{O(1)}$ translates of some subspaces of size at most $|S|$.*

By the probabilistic bound in (10), we are therefore motivated to state the following conjecture.

Conjecture 2. *If the linear compression dimension of S is given by $\log_2(|S|) + o(n)$ (hence matches the non-linear compression dimension) then S is contained in the union of $2^{o(n)}$ translates of some subspaces of size at most $|S|$.*

The above condition can only happen for sets which are non symmetric, i.e., not invariant under permutations of the n coordinates (like a subspace), unless the set is very small (of size $2^{o(n)}$) or very large (of size $2^{n+o(n)}$). In the next section, we will focus on symmetric sets.

C. Linear Boolean classification for general models

Definition 3. *Let $S_1, S_2 \subseteq \{0, 1\}^n$ with $S_1 \cap S_2 = \emptyset$. A linear classifier for (S_1, S_2) is a linear map $M : \mathbb{F}_2^n \rightarrow \mathbb{F}_2^m$ which is (S_1, S_2) -separable in the sense that*

$$Mx_1 \neq Mx_2, \quad \forall x_1 \in S_1, x_2 \in S_2. \quad (12)$$

The linear classification dimension of the set pair (S_1, S_2) is defined by

$$m^*(S_1, S_2) = \min_{M \in (S_1, S_2)\text{-separable}} \text{rank}(M). \quad (13)$$

Note that $m^*(S_1, S_2)$ can also be expressed as

$$m^*(S_1, S_2) = n - \max_{V: V \cap (S_1 + S_2) = \emptyset} \dim(V), \quad (14)$$

where V denotes a subspace of \mathbb{F}_2^n . In other words, we need to find the largest dimension of a subspace avoiding $S_1 + S_2$. This is again an instance of “the critical problem”. The linear classification rate is defined by $m^*(S_1, S_2)/n$.

Note however that the two simple bounds obtained in the previous section do not give interesting results here. First, if non-linear maps are allowed, the lower bound on m^* is simply 1 as there are only two sets. In terms of probabilistic bounds, if M is drawn uniformly at random, then we obtain

$$m^*(S_1, S_2) \leq \lceil \log_2(|S_1 + S_2|) \rceil. \quad (15)$$

Unlike in the coding problem, this is not a strong bound in general, as shown in Section IV. In particular, it does

not take advantage of the fact that $S_1 + S_2$ is away from 0^n .

III. DEFINITIONS OVERVIEW AND PROBLEM STATEMENT

Let $n \in \mathbb{Z}_+$ and $S, S_1, S_2 \subseteq \{0, 1\}^n$.

Definition 4. A matrix M is

- S -distinguishable if $Mx \neq 0$ for all $x \in S$,
- S -injective if $Mx \neq Mx'$ for all $x, x' \in S, x \neq x'$,
- (S_1, S_2) -separable if $Mx_1 \neq Mx_2$ for all $x_1 \in S_1, x_2 \in S_2$.

Note that

- M is S -distinguishable iff $\ker(M) \cap S = \emptyset$,
- M is S -injective iff M is $((S + S) \setminus \{0\})$ -distinguishable,
- M is (S_1, S_2) -separable iff M is $(S_1 + S_2)$ -distinguishable.

From the first item above, a matrix of minimal rank which is S -distinguishable can be equivalently constructed by finding a space of maximal dimension which avoids S , an instance of “the critical problem”. The meaning of previous mathematical objects in terms of compression, coding and classification notions are:

- M is a lossless compressor for the source model S iff M is S -injective,
- M is the parity-check matrix of an error-correcting code for the error model S iff it is a lossless compressor for the source model S ,
- M is a linear classifier for the class models (S_1, S_2) iff M is (S_1, S_2) -separable.

One could also use classification in a coding context for errors, to determine if an error pattern belongs to two different classes (e.g., typical or atypical patterns, high or low SNR regimes). In this case the linear constraint on the classifier is important to allow the source-channel coding duality.

In this paper, we are interested in symmetric sets, i.e., sets which are invariant under permutations. Note that these are defined by the Hamming weights of their elements. One of the most fundamental symmetric sets is the Hamming ball at 0^n , studied extensively in coding theory. The only other symmetric Hamming ball is the one centered at 1^n . A natural set structure to consider next is the annulus, which contains the cases above and extends to more general symmetric sets. In particular, an arbitrary symmetric set is a union of annuli. In addition, the sum of two annuli, which matters for classification, is again an annulus (except in particular cases where the two annuli have a single weight, and their sum contains only even or odd weight vectors).

Definition 5. Let $1 \leq a \leq b \leq n$,

$$A(a, b, n) = \{x \in \{0, 1\}^n : w(x) \in [a, b]\}, \quad (16)$$

$$m^*(a, b, n) = \min_{M \in A(a, b, n)\text{-distinguishable}} \text{rank}(M) \quad (17)$$

$$= n - \max_{V: V \cap A(a, b, n) = \emptyset} \dim(V), \quad (18)$$

where M is a matrix over \mathbb{F}_2 with n columns and V is a subspace of \mathbb{F}_2^n .

Note that $m^*(1, b, n) = m^*(b, n)$, i.e., when the annulus degenerates into the punctured ball $B(0^n, b) \setminus \{0^n\}$, the definition is consistent with the one of Section I-A.

Our goal is to characterize $m^*(a, b, n)$ for finite values of the parameters and in the asymptotic regime.

IV. RESULTS

Recall that $m^*(1, b, n)$, the least rank of a parity-check matrix of a code of distance $b + 1$ on a blocklength n , has no known explicit form in general. It is clear that $m^*(1, b, n) = n$, if $b = n$, and it is known that

$$m^*(1, b, n) = n + o(n), \quad \text{if } b \geq n/2, \quad (19)$$

from Plotkin's bound [8]. Also, the GV bound provides the inequality

$$m^*(1, \beta n, n) \leq nH(\beta) + o(n), \quad (20)$$

which is conjectured to be tight. In this paper we are interested in characterizing $m^*(a, b, n)$ for $a \geq 2$, possibly in terms of $m^*(1, b, n)$. In other words, we want to study how the "hole" in an annulus allows one to decrease the dimension of M when compared to a Hamming ball.

We obtain the following result.

Proposition 1. For any $1 \leq a \leq b \leq n$ with $b \geq 2a - 2$,

$$m^*(a, b, n) = m^*(1, b, n - a + 1). \quad (21)$$

Assuming the GV bound to be tight, the above result takes the following asymptotic form

$$m^*(\alpha n, \beta n, n) = n(1 - \alpha)H\left(\frac{\beta}{1 - \alpha}\right) + o(n). \quad (22)$$

In the proof of this proposition, we show that an optimal kernel for an $A(a, b, n)$ -distinguishable matrix is composed by a mixture of $a - 1$ unit vectors and a code of largest dimension and minimal distance at least $b + 1$ on the remaining components.

Using (19), we obtain the following corollary which provides a characterization of m^* in certain cases.

Corollary 1. If $b \geq 2a - 2$ and $2b \geq n - a + 1$,

$$m^*(a, b, n) = n - a + 1 + o(n). \quad (23)$$

It is also straightforward to show that the following holds for any a :

$$m^*(a, n, n) = n - a + 1. \quad (24)$$

The proof simply uses the fact that a subspace containing a linearly independent vectors must generate a vector of weight at least a .

Going back to the toy example of Section I-B, where $S_1 = B(1^n, s)$ and $S_2 = B(0^n, s)$, $s < n/2$, are the two classes, we have $S_1 + S_2 = B(1^n, 2s) = A(n - 2s, n, n)$. Hence, by (24), $m^* = 2s + 1$. This is indeed verified by the fact that the identity matrix of dimension $2s + 1$ allows to classify the two sets (by a simple majority count on any $2s + 1$ coordinates). This also shows that the classification dimension can be simpler to find than the compression dimension, since the sum of two disjoint sets $S_1 + S_2$ is typically away from 0^n and hence V can contain sparse vectors, whereas the sumset $S + S$ contains a ball around 0^n and requires a packing of dense vectors, which is more challenging.

We conjecture that Proposition 1 holds unconditionally, besides for degenerated cases of the form

$$m^*(a, a, n) = 1, \quad a \text{ odd, } n \text{ even.} \quad (25)$$

It is straightforward to verify (25) by noting that packing vectors of even weight will never produce an odd weight vector.

Conjecture 3. Proposition 1 holds for all $1 \leq a < b \leq n$, except for degenerate cases of the form of (25).

The upper-bound holds in full generality.

Proposition 2. For any $1 \leq a \leq b \leq n$,

$$m^*(a, b, n) \leq m^*(1, b, n - a + 1). \quad (26)$$

To achieve the above, it is enough to take for the kernel of M a basis consisting of $a - 1$ unit vectors and an optimal subspace of weight at least $b + 1$ on the complement coordinates.

We further obtain a few more cases which corroborate the conjecture.

In particular, as observed by the second author in the 80s (c.f. [4]), a theorem of Olson [7] can be used to show that if a is a power of 2 then any binary linear code of dimension a (and any length) contains a vector of Hamming weight divisible by a . This implies that if $a > n/2$ is a power of 2 then $k(a, a, n) = n - m^*(a, a, n) = a - 1$. A more general result is proved in [4] where it is shown that for any even $a > n/2$, $k(a, a, n) = a - 1$. Further results on codes with a forbidden distance can be found in [2].

V. FUTURE WORK

For symmetric sets, the next steps would be to investigate further regimes for the radius of the annulus, supporting or disproving the conjecture or to consider the union of two annuli. It would also be interesting to consider a probabilistic rather than worst-case model for the linear Boolean classification problem. The complexity of the classification would be another interesting direction to pursue. Finally, a natural extension is to consider the problem of constructing matrices that allow to classify certain sets while compressing others. This will be another instance of “the critical problem”.

VI. PROOFS

In the proofs we work with the kernel approach, maximizing the dimension of $V = \ker(M)$ without intersecting the annulus $A(a, b, n)$. We use the notation $k(a, b, n) = n - m^*(a, b, n)$ and provide a proof of Proposition 1 based on the weight of the sparsest vector in an optimal basis.

Lemma 2. *If $b \geq 2a - 2$ and $a \geq 2$,*

$$k(a, b, n) \leq \max_{1 \leq s \leq a-1} [s + k(a - s, b, n - s)].$$

Proof of Lemma 2:

Let V be a subspace that does not intersect $A(a, b, n)$, and s be the sparsity of the sparsest non-zero vector in V . If $s \geq a$ then $\dim(V) \leq k(1, b, n)$ and using $k(1, b, n - 1) + 1 \geq k(1, b, n)$,

$$\dim(V) \leq k(1, b, n) \leq 1 + k(1, b, n - 1) \quad (27)$$

$$\leq 1 + k(a - s, b, n - 1) \quad (28)$$

$$\leq \max_{1 \leq s \leq a-1} [s + k(a - s, b, n - s)]. \quad (29)$$

On the other hand, if $s < a$, we will show that $\dim(V) \leq s + k(a - s, b, n - s)$, proving the Lemma.

Let v be a vector in V that is exactly s -sparse. We permute (the coordinates of) V so that v is 1 in the first s components and 0 elsewhere. We represent V by a matrix (below) with its rows forming a basis of V , we pick such a representation such that v is the first row.

$$\begin{bmatrix} 1_{1 \times s} & 0_{1 \times (n-s)} \\ \vdots & \vdots \end{bmatrix},$$

If $\dim(V) < s$ then the result is trivial, so we will focus on the case $\dim(V) \geq s$, using Gauss Elimination one can find a basis of V such that v is the first row and the first s by s block is upper triangular, meaning:

$$\begin{bmatrix} 1 & 1_{1 \times (s-1)} & 0_{1 \times (n-s)} \\ 0_{(s-1) \times 1} & T & R \\ 0_{(\dim(V)-s) \times 1} & 0_{(\dim(V)-s) \times (s-1)} & V^* \end{bmatrix},$$

with T upper-triangular.

Note that V^* is a basis for a subspace in $n - s$ coordinates of dimension $\dim(V) - s$. We next argue that $V^* \cap A(a - s, b, n - s) = \emptyset$. Indeed, suppose $u \in V^* \cap A(a - s, b, n - s)$ then

- If $a \leq w(u) \leq b$ then the vector $[0_{s \times 1} \ u] \in V$ is in $A(a, b, n)$.
- On the other hand, if $a - s \leq w(u) < a$ the vector $[0_{s \times 1} \ u] \in V$ summed to the s sparse vector $v \in V$ will give a vector in $A(a, a - 1 + s, n)$. Since $a - 1 + s \leq a - 1 + a - 1 = 2a - 2 \leq b$ then $A(a, a - 1 + s, n) \subset A(a, b, n)$.

This means that $\dim(V^*) \leq k(a - s, b, n - s)$ and $\dim(V) - s \leq k(a - s, b, n - s)$. ■

The proof of Proposition 1 follows then by a strong induction on a and Proposition 2.

Remark 1. Note that the proof above would carry through if $b < 2a - 2$, as long as there exists an element in the optimal subspace, whose Hamming weight is smaller or equal to $b - a + 1$. This means that, for Conjecture 3, the subspace must have all its elements' weights not only avoiding $[a, b]$ but also avoiding $[1, b - a + 1]$.

Remark 2. An intuitive way of thinking about the condition $b \geq 2a - 2$ is that it enforces that any sum of sparse vectors $x, y \in A(a, b, n)$ (meaning $w(x), w(y) \leq a - 1$) cannot be dense, as $w(x + y) \leq b$. Indeed, one can use this fact to provide an alternative proof to Proposition 1.

REFERENCES

- [1] E. Abbe, *Worst-case source coding*, Course notes: Coding Theory and Random Graphs, Princeton University. Available at www.princeton.edu/eabbe, 2013.
- [2] L. Bassalygo, G. Cohen, and G. Zémor, *Codes with forbidden distances*, Discrete Mathematics, Volume 213, Issues 1–3, Pages 3–11, February 2000.
- [3] H.H. Crapo and G.C. Rota, *On the foundations of combinatorial theory: Combinatorial geometries*, MIT Press, Cambridge, MA, 1970.
- [4] H. Enomoto, P. Frankl, N. Ito, and K. Nomura, *Codes with given distances*, Graphs and Combinatorics **3** (1987), no. 1, 25–38.
- [5] E.N. Gilbert, *A comparison of signalling alphabets*, Bell System Technical Journal **31** (1952), no. 3, 504–522.
- [6] B. Green, *Notes on the polynomial Freiman-Ruzsa conjecture*, available online (2005), <http://people.maths.ox.ac.uk/greenbj/papers/PFR.pdf>.
- [7] J.E. Olson, *A combinatorial problem on finite abelian groups*, i, Journal of Number Theory **1** (1969), no. 1, 8–10.
- [8] M. Plotkin, *Binary codes with specified minimum distance*, Information Theory, IRE Transactions on **6** (1960), no. 4, 445–450.
- [9] R.R. Varshamov, *Estimate of the number of signals in error correcting codes*, Dokl. Acad. Nauk SSSR **117** (1957), 739–741.