

On the Mathematical Theory of Error-Correcting Codes

Abstract: Hamming considered the problem of efficient, faultless transmission of binary data over a noisy channel. For a channel which corrupts no more than one binary digit in each sequence of length n , he constructed alphabets, the so-called Hamming codes, which permit error-free signalling. The authors study the analogous problem for channels which can corrupt a greater number of digits. Non-binary channels are also studied, and analogues of the Hamming codes are constructed. It is perhaps of interest that some of the techniques employed derive from algebraic and analytic number theory, mathematical disciplines not generally associated with the type of applied problems considered in this paper.

1. Introduction and survey of results

We must begin by introducing sufficient notation to discuss the background of the subject and describe our results. To this end let \mathcal{G}_n denote the set of binary sequences of length n . There are $N=2^n$ such sequences or "points" which may be considered as the coordinates of the vertices of a unit cube in n -space. A set \mathcal{C} consisting of M such vertices is called, following Hamming,³ an e -error correcting code if any two points of \mathcal{C} differ in at least $2e+1$ coordinates. The reason for this designation is that if e or fewer "errors" are made in transmitting a binary sequence of \mathcal{C} , i.e., if not more than e digits undergo change during transmission, it is possible in principle to recover the transmitted sequence. Designating by the distance between two points x, y of \mathcal{G}_n the number of coordinates in which they differ, we are thus led to the problem of finding subsets \mathcal{C} of \mathcal{G}_n whose points have mutual distance $\geq d=2e+1$; such a set is called an (n, d) code.

Following the notation of Plotkin,⁸ we let $A(n, d)$ denote the maximum number of points of \mathcal{G}_n which can have mutual distance $\geq d$. This definition is, of course, meaningful also for d even, $d=2e$. A set of vertices having mutual distance $\geq 2e$ is called an e -error detecting code because in this case it is possible to completely restore a sequence containing $\leq e-1$ changes and to recognize, given a sequence containing e changes, that an error in transmission has been made, although it is in general not possible to restore the message unambiguously. It is readily shown that all problems concerning e error detecting codes are reducible to problems concerning $(e-1)$ -error-correcting codes of order $n-1$, i.e., $A(n, 2e)=A(n-1, 2e-1)$, and a simple correspondence exists between $(n, 2e)$ and $(n-1,$

$2e-1)$ codes involving a single parity check; thus error-detecting codes as such will not be further considered in this paper.

Under digit-wise modulo 2 addition \mathcal{G}_n is a group, and the distance between two points x, y is just the number of ones in the sequence $x+y$; we denote by $\|x\|$ (read: *norm* or *weight* of x) the number of ones in the sequence x . Of special interest are those (n, d) codes \mathcal{C} which are subgroups of \mathcal{G}_n ; such codes are called *systematic* or *group codes*, and are of special interest for the following reason: the decoding of received messages, i.e., the actual correction of wrong digits, can be effected by means of a *parity-check matrix* if and only if the code employed is a group code.* Another advantage is that group codes can be completely specified by giving relatively few data, namely a basis or even a parity-check matrix. Although it is likely that certain non-group codes also lend themselves to simple mathematical specification and detection procedures other than parity checks,[†] the group codes are of central importance. They are also especially convenient to work with mathematically. Because $x+y$ is a code point whenever x and y are, we see that a *systematic (group) e -error correcting code is a subgroup of \mathcal{G}_n all of whose elements (except $0=00\dots 0$) contain at least $d=2e+1$ ones*. The maximum size of such a code we denote by $B(n, d)$ and a code attaining this size is called a *maximal* code. The remarks about error-detecting codes apply also to group codes, i.e., we have

*For a proof, and a very lucid discussion of parity-check procedures and other basic matters, see Slepian⁹. See also Lloyd⁶.

†In this category are the "quadratic residues" codes of Plotkin⁸, which one can generalize on the basis of work of Paley⁷ and others in connection with the quite different problem of "Hadamard matrices".

$B(n, 2e) = B(n-1, 2e-1)$ and a simple correspondence between respective codes of the two kinds. Obviously $B(n, d)$ is a power of 2 and $B(n, d) \leq A(n, d) \leq 2^n$. We have further, as Hamming³ showed,

$$A(n, 2e+1) \leq \frac{2^n}{1 + \binom{n}{1} + \dots + \binom{n}{e}}. \quad (1)$$

To see this, associate with each point x of an $(n, 2e+1)$ code \mathcal{C} the "sphere" \mathcal{S}_x consisting of all points of \mathcal{G}_n having distance $\leq e$ from x . If $x \neq y$, \mathcal{S}_x and \mathcal{S}_y are disjoint; for if not, a point z lies in both \mathcal{S}_x and \mathcal{S}_y , i.e., has distance $\leq e$ from x and y , whence x and y have distance $\leq 2e$, a contradiction. Hence, the number of code points is precisely the number of spheres, and since each \mathcal{S}_x contains $1 + \binom{n}{1} + \dots + \binom{n}{e}$ points, and \mathcal{G}_n contains 2^n points altogether, the number of spheres cannot exceed the right side of (1).

In the case $e=1$, we have

$$A(n, 3) \leq \frac{2^n}{n+1}. \quad (2)$$

The right side is an integer if and only if n has the form $2^k - 1$ and Hamming showed that in this case equality holds in (2), in fact even

$$B(n, 3) = \frac{2^n}{n+1}, \quad (3)$$

that is, equality is attained for a *group* code. And for all n , whether or not of the form $2^k - 1$, $B(n, 3)$ is the largest power of 2 not exceeding $2^n/(n+1)$. Any code which gives equality in (2) corresponds to a remarkable decomposition of \mathcal{G}_n , for in this case the disjoint spheres \mathcal{S}_x described above completely exhaust \mathcal{G}_n . Hence, such codes, whether systematic or not, are called *close-packed*. Close-packed codes, when they exist, are maximal and also have many remarkable symmetries, as we shall see below. It is natural to begin a general study of error-correcting codes with a quest for close-packed codes. We will see that the results for $e \geq 2$, unlike the case $e=1$, are basically negative, i.e., there are no close-packed double-error correcting codes except the trivial (5, 5) code with 2 points; for $e=3$, there is a remarkable group code that was discovered by Golay, a close-packed (23, 7) code having 2^{11} points. We will show that, aside from this code and the trivial (7, 7) code having two points, there are no further close-packed triple-error correcting codes. Detailed settlement of the question for $e \geq 4$ is more difficult, although we can show easily that for each $e \geq 4$ the number of close-packed e -error correcting codes is finite. Of course, this does not preclude the existence of very large and non-trivial close-packed codes. For e odd, moreover, we shall give a procedure for obtaining any possible close-packed code.

In Section 4 we obtain a new lower bound for $B(n, d)$. The gap between this bound and Hamming's upper bound is still great enough to have qualitative significance, in the sense of attaining channel capacity⁸ for a symmetric binary channel, by means of error-correcting codes. We give also a procedure by means of which large (n, d) systematic codes can be constructed for successive values of n ; although

some trial and error is involved, the procedure lends itself easily to machine computation.

In Section 5 we extend the Hamming codes and the results of Section 4 to the case of a p -ary channel, where p is any prime number.* Although the p -ary channel is presently of less practical interest than the binary, fairly large binary codes can be constructed by means of p -ary codes, by encoding p -ary symbols into binary.

2. Symmetries of close-packed codes

• Theorem 1

Let \mathcal{C} be a close-packed (n, d) code ($d=2e+1$). Then the ratio

$$\nu = \binom{n}{e+1} / \binom{2e+1}{e}$$

is an integer, and every code point has distance d from precisely ν other code points.†

Proof

Let x be a point of \mathcal{C} . Without loss of generality we may assume that $x=0$. For, if we replace each point y of \mathcal{C} by $y+x$, we get a configuration congruent to \mathcal{C} containing $0=x+x$. Consider now the set \mathcal{Y} of all points y of \mathcal{G}_n at distance $e+1$ from 0 , i.e., points of \mathcal{G}_n having exactly $e+1$ ones. There are $\binom{n}{e+1}$ such. Each y lies in a sphere of radius e about some particular code point z , since \mathcal{C} is close-packed. Clearly $\|z\|=2e+1$. Thus, each such y is obtained once and only once by taking a code point having $2e+1$ ones, and changing e of its ones to zeros, i.e., the number of points in \mathcal{Y} is $\binom{2e+1}{e}$ times ν , the number of code points z of norm $2e+1$. This proves the theorem.

• Theorem 2

Let \mathcal{C} be a close-packed $(n, 2e+1)$ code, and x and y points of \mathcal{C} . Then, given any $e+1$ distinct integers of the set $1, 2, \dots, n$ there is one and only one point y of \mathcal{C} having distance $2e+1$ from x , and differing from x in each of these $e+1$ positions.

Proof

As in Theorem 1, it is sufficient to assume $x=0$. Now, that there is at most one y with the stated properties is clear; for, if there were two such, y_1 and y_2 , we would have $\|y_1\| = \|y_2\| = 2e+1$ and, since y_1 and y_2 must agree in at least $e+1$ digits (namely, those digits where they differ from $x=0$, i.e., where they have ones), we should have $\|y_1 + y_2\| \leq 2e$, a contradiction. On the other hand, let $f(I)$ denote the number of code points of norm $2e+1$, having ones in $e+1$ prescribed positions (here I denotes a set of $e+1$ distinct positions). Then,

*Added in proof: The authors have since learned of work similar to that in Section 5 by Golay⁴, Lee⁵, and Ulrich.¹¹

†A system of equations apparently equivalent to our system (S) derived in Section 2, but involving complicated sums of binomial coefficients was obtained by Lloyd⁶ using a generating function technique.

$$\sum f(I) = \binom{2e+1}{e+1} \cdot \nu = \binom{n}{e+1},$$

since each code point of norm $2e+1$ is counted $\binom{2e+1}{e+1}$ times in this summation. However, since $f(I) = 0$ or 1 (by the above) and there are precisely $\binom{n}{e+1}$ summands on the left-hand side, $f(I)$ is always 1 , otherwise the sum on the left would be $< \binom{n}{e+1}$, a contradiction.

• **Theorem 3**

Let \mathcal{C} be a close-packed $(n, 2e+1)$ code. Then $\mu = \frac{n-e}{e+1}$ is an integer, and if \mathbf{x} is any code point, the number of \mathbf{y} in \mathcal{C} at distance $2e+1$ from \mathbf{x} , and differing from \mathbf{x} in e prescribed positions is precisely μ .

Proof

As before, we may suppose $\mathbf{x} = \mathbf{0}$. Let I denote a set of e distinct positions, and $f(I) =$ the number of code points of norm $2e+1$ which differ from \mathbf{x} in all of the positions I . As before

$$\sum f(I) = \binom{2e+1}{e} \cdot \nu.$$

On the other hand,

$$f(I) \leq \left\lceil \frac{n-e}{e+1} \right\rceil$$

for every I , where $\lceil \cdot \rceil$ denotes "greatest integer", for, consider the set of \mathbf{y} of norm $2e+1$ having ones in e places; the remaining $e+1$ ones must be placed in mutually disjoint positions or else minimum distance $2e+1$ cannot be

preserved. Thus, there are at most $\left\lceil \frac{n-e}{e+1} \right\rceil$ such \mathbf{y} . Hence,

since the number of admissible I is $\binom{n}{e}$, we have

$$\binom{n}{e} \left\lceil \frac{n-e}{e+1} \right\rceil \geq \binom{2e+1}{e} \cdot \nu = \binom{n}{e+1}.$$

Hence,

$$\binom{n}{e} \left\lceil \frac{n-e}{e+1} \right\rceil \geq \binom{n}{e} \cdot \frac{n-e}{e+1},$$

and equality can hold only if

$$\left\lceil \frac{n-e}{e+1} \right\rceil = \frac{n-e}{e+1},$$

and further $f(I) = \frac{n-e}{e+1}$ for all I . This completes the proof.

• **Theorem 4**

Let \mathcal{C} be a close-packed $(n, 2e+1)$ code and $2e+1 \leq r \leq n$. Then the number of code points having distance precisely r from a given code point \mathbf{x} is an integer $\nu(n, e, r)$ which does not depend on the particular point \mathbf{x} nor upon the particular code chosen.

Proof

For convenience of notation we set $\nu(n, e, r) = a_r$. The proof will obtain by deriving a system of linear equations for the a_r from which they can actually be calculated. For simplicity of presentation we deal with the concrete case $e=2$, but the procedure employed is perfectly general.

Consider the subset \mathcal{G}_{nk} of \mathcal{G}_n , consisting of those elements of \mathcal{G}_n of norm k . Since \mathcal{C} is close-packed, we may with every \mathbf{x} in \mathcal{G}_n associate a unique $\mathbf{y} = \mathbf{y}(\mathbf{x})$ in \mathcal{C} such that $\|\mathbf{x} + \mathbf{y}\| \leq 2$. In particular we may do this for \mathbf{x} in \mathcal{G}_{nk} and thus decompose \mathcal{G}_{nk} into three mutually disjoint subsets

$$\mathcal{G}_{nk} = \mathcal{G}_{nk}^0 + \mathcal{G}_{nk}^1 + \mathcal{G}_{nk}^2,$$

where an \mathbf{x} of \mathcal{G}_{nk} is placed in \mathcal{G}_{nk}^r if $\|\mathbf{x} + \mathbf{y}(\mathbf{x})\| = r$, ($r=0, 1, 2$). \mathcal{G}_{nk}^0 consists of those \mathbf{x} in \mathcal{G}_{nk} which belong to \mathcal{C} , and so has a_k elements. If \mathbf{x} is in \mathcal{G}_{nk}^1 , we may distinguish two cases:

- (i) $\|\mathbf{y}(\mathbf{x})\| = k-1$
- (ii) $\|\mathbf{y}(\mathbf{x})\| = k+1$.

The number of \mathbf{x} in \mathcal{G}_{nk}^1 for which (i) holds is $a_{k-1}(n-k+1)$, i.e., the number of ways we can pick a code point of norm $k-1$ and then change one of its zero-digits to a 1. The number of \mathbf{x} in \mathcal{G}_{nk}^1 for which (ii) holds is $a_{k+1}(k+1)$, i.e., the number of ways in which we can pick a code point of rank $k+1$ and change one of its ones to a zero.

In like manner we may partition the \mathbf{x} in \mathcal{G}_{nk}^2 into three classes according as $\|\mathbf{y}(\mathbf{x})\| = k-2$, k , or $k+2$ and these classes contain

$$\binom{n-k+2}{2} a_{k-2}, k(n-k) a_k, \binom{k+2}{2} a_{k+2}$$

elements respectively. For instance, the second of these numbers is the number of ways we can choose a code point of norm k , then change simultaneously one of its ones to a zero and one of its zeroes to a one.

Now, since \mathcal{G}_{nk} contains $\binom{n}{k}$ elements, we deduce from the above considerations that*

$$\begin{aligned} \binom{n}{k} = & a_k + \left[(n-k+1) a_{k-1} + (k+1) a_{k+1} \right] \\ & + \left[\binom{n-k+2}{2} a_{k-2} + k(n-k) a_k + \binom{k+2}{2} a_{k+2} \right] \end{aligned}$$

for $k=0, 1, \dots, n$.

This is a system of $n+1$ linear equations for the $n+1$ "unknowns" a_0, a_1, \dots, a_n .

If, for a given n , a close-packed $(n, 5)$ code exists, these equations are a singular system, i.e., the determinant of the system must vanish. For otherwise the values a are determined independently of the particular code, in particular a_0 is determined, so that either all codes contain the origin or all do not. But this is impossible since as before we can go from a code containing the origin to a congruent one not containing the origin (and vice versa) by a simple translation.

*Negatively indexed a are to be interpreted as zero.

On the other hand, supposing that \mathfrak{S} contains the origin, i.e., that $a_0=1$ (whence $a_1=a_2=a_3=a_4=0$) reduces the system to one which can be solved for $a_5, a_6, a_7, \dots, a_n$, viz:

$$\begin{aligned} 10a_5 &= \binom{n}{3} \\ 5a_5 + 15a_6 &= \binom{n}{4} \\ &\text{etc.} \quad \dots \end{aligned}$$

Similarly, in the general case we get (assuming $a_0=1$)
 $a_2=a_3=\dots=a_{2e}=0$,

$$\begin{aligned} \binom{2e+1}{e} a_{2e+1} &= \binom{n}{e+1} \\ \binom{2e+1}{e-1} a_{2e+1} + \binom{2e+2}{e} a_{2e+2} &= \binom{n}{e+2} \\ \binom{2e+1}{e-2} a_{2e+1} + \binom{2e+2}{e-1} a_{2e+2} + \binom{2e+3}{e} a_{2e+3} &= \binom{n}{e+3} \\ &\vdots \end{aligned} \quad (S)$$

The law of formation of these equations is evident, from which the a_r can be successively calculated. This proves Theorem 4.

Remark 1

A necessary condition that there exist a close-packed (n, d) code is thus that the a_r obtained by solving this system (S) be integers.

Remark 2

Theorem 4 suggests, although it does not follow directly from the theorem, that there is essentially at most one close-packed (n, d) code. In other words, given two close-packed (n, d) codes there is a symmetry of the n -cube which takes one into the other.

In the case of the Hamming codes it is not hard to solve for the a_r explicitly (e.g., by introducing the generating function $a_0 + \dots + a_n t^n$); we get

$$\begin{aligned} (n+1)a_{2r} &= \binom{n}{2r} + (-1)^r n \binom{\frac{n-1}{2}}{r} \\ (n+1)a_{2r+1} &= \binom{n}{2r+1} + (-1)^{r+1} n \binom{\frac{n-1}{2}}{r} \end{aligned}$$

These equations display the finer structure of the Hamming codes; an interesting consequence is

$$a_{2r} + a_{2r+1} = \frac{\binom{n}{2r} + \binom{n}{2r+1}}{n+1}.$$

This equation states that if we consider successive "layers" of \mathfrak{G}_n consisting first of points of rank 0 or 1, then points of rank 2 or 3, then points of rank 4 or 5, etc., the code points are, so to speak, uniformly distributed within each layer.

We are now ready to consider whether close-packed error-correcting codes of higher order can exist.

3. The existence problem for close-packed codes,* $e \geq 2$.

• Theorem 5

The only close-packed double-error correcting code is the trivial $(5, 5)$ code with 2 points.

The proof of this is surprisingly difficult and employs the arithmetic of the algebraic number field $\mathfrak{f}(\sqrt{-7})$, about which we shall require the following information (see, e.g., Hecke[†]; the reader may skip this proof if he so desires, without prejudice to the remaining theorems).

Lemma

In the algebraic number field $\mathfrak{f}(\sqrt{-7})$ all integers are of the form $a+b\rho$, where a, b are rational integers and $\rho = \frac{-1+\sqrt{-7}}{2}$ is a root of $x^2+x+2=0$. The only units are ± 1 , and factorization is unique. Moreover,

$\rho, \bar{\rho} = \frac{-1-\sqrt{-7}}{2}$, and $\sqrt{-7} = \rho - \bar{\rho}$ are primes of $\mathfrak{f}(\sqrt{-7})$.

Suppose now that a close-packed $(n, 5)$ code exists. Then we must have (i) $1+n+\binom{n}{2}$ divides 2^n , since for equality to hold in Hamming's inequality (1) the right side must be an integer, and also (ii) $n \equiv 2 \pmod{3}$ by virtue of Theorem 3. We will show that for $n > 5$, (i) and (ii) cannot hold simultaneously. Since any divisor of a power of 2 is itself a power of 2, we may write

$$1+n+\binom{n}{2} = 2^k, \text{ or} \quad (4)$$

$$n^2+n+2 = 2^{k+1}. \quad (5)$$

Let us now consider Eq. (5) as an equation between integers of $\mathfrak{f}(\sqrt{-7})$, which we may certainly do since the rational integers are integers of $\mathfrak{f}(\sqrt{-7})$. We have

$$(n-\rho)(n-\bar{\rho}) = \rho^{k+1}\bar{\rho}^{k+1}. \quad (6)$$

Now, the factors on the left differ by $\rho - \bar{\rho} = \sqrt{-7}$ which is not divisible by ρ or $\bar{\rho}$, hence in view of the lemma one of the following four cases holds:

- (a) $n-\rho = \rho^{k+1}$
- (b) $n-\rho = -\rho^{k+1}$
- (c) $n-\rho = \bar{\rho}^{k+1}$
- (d) $n-\rho = -\bar{\rho}^{k+1}$,

and we must show that none of these cases holds for $n > 5$, simultaneously with $n \equiv 2 \pmod{3}$. In fact, (b), (c), (d) can be excluded without recourse to this latter condition. We eliminate first (b) and (c): in each case we have

$$a_k = \frac{\rho^{k+1} - \bar{\rho}^{k+1}}{\rho - \bar{\rho}} = 1,^\dagger \quad (7)$$

Now, the a_k are rational integers[†] and

*The referee has pointed out that Golay² has considered this question in the case of parity-check codes, and obtained our Theorems 5 and 8 for this special class of codes. The analysis is greatly simplified in the case of parity-check codes by the presence of an additional condition [Golay's Eq. (2)].

†These a_k , of course, have nothing to do with the a_k of Section 2.

once again, we deduce, since $k \geq 3$ by assumption

$\rho|r+1$ i.e., $r+1$ is even $= 2s$.

Hence $n = 2m = 4l - 2 = 8r + 2 = 16s - 6$

and so

$$n \equiv -6 \pmod{16}. \quad (24)$$

Thus, for case (a) to hold with $k \geq 3$ (and $n \equiv 2 \pmod{3}$) we must have, from (19) and (24)

$$a_{k+1} \equiv -5 \pmod{16}. \quad (25)$$

But the residues of the $a_k \pmod{16}$ are as follows:

$$\begin{array}{c|c|c|c|c|c|c|c|c|c|c|c} k & 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 \\ \hline a_k \pmod{16} & 1 & -1 & -1 & 3 & -1 & -5 & 7 & 3 & -1 & -5 & 7 \end{array}.$$

That is, they are 1, -1, -1 followed by the repeated period 3, -1, -5, 7. Hence (25) holds only if

$$k+1 \equiv 5 \pmod{4}, \text{ i.e.,}$$

$$k \equiv 0 \pmod{4} \quad (26)$$

(26) and (21) stand in contradiction, and the proof is completed.

• Theorem 6

The only close-packed triple-error correcting codes are the (23, 7) code of Golay¹ and the trivial (7, 7) code.

Proof

$$\text{If } 1 + \binom{n}{1} + \binom{n}{2} + \binom{n}{3} = 2^k$$

then, on multiplying by 6 and simplifying the left side factors, we get

$$(n^2 - n + 6)(n + 1) = 3 \cdot 2^{k+1}. \quad (27)$$

Hence, one or the other of the left-hand factors is a multiple of 3.

$$\text{Case I.} \quad 3|n^2 - n + 6.$$

$$\text{Here } n+1 = 2^l, \quad n^2 - n + 6 = 3 \cdot 2^{k+1-l}$$

$$\text{whence, } (2^l - 1)^2 - (2^l - 1) + 6 = 3 \cdot 2^{k+1-l}$$

$$2^{2l} - 3 \cdot 2^l + 8 = 3 \cdot 2^{k+1-l}. \quad (28)$$

Obviously we may restrict ourselves to the case $n > 7$ (for $n = 7$, we get the trivial code). Hence, $l \geq 4$.

Since the left side of (28) is $\equiv 8 \pmod{16}$ so is the right, and thus $k+1-l=3$ and $n^2 - n + 6 = 24$, which contradicts $n > 7$ (as well as $n = \text{rational integer}$).

$$\text{Case II.} \quad 3|n+1$$

$$\text{Here, } n+1 = 3 \cdot 2^l \text{ and } n^2 - n + 6 = 2^{k+1-l},$$

$$(3 \cdot 2^l - 1)^2 - (3 \cdot 2^l - 1) + 6 = 2^{k+1-l}$$

$$9 \cdot 2^{2l} - 9 \cdot 2^l + 8 = 2^{k+1-l} \quad (29)$$

$$\text{since } n \geq 8, \quad l \geq 2 \text{ and } k+1-l \geq 5.$$

$$\text{From (29), } 9 \cdot 2^l \equiv 8 \pmod{16}.$$

Hence, $l = 3$, $n+1 = 3 \cdot 2^3$, $n = 23$, completing the proof; this

leads to the (23, 7) code mentioned in the theorem. If it is desired to construct the code without simply quoting Golay's construction, we outline briefly how it could be done (with relatively little trial and error). The trick is to start by looking for code points of norm 7. By Theorem 3, there is no loss of generality in assuming as code points (taking also $\mathbf{x} = \mathbf{0}$ as a code point).

$$\mathbf{x}_1 = 111 \quad 1111 \quad 0000 \quad 0000 \quad 0000 \quad 0000$$

$$\mathbf{x}_2 = 111 \quad 0000 \quad 1111 \quad 0000 \quad 0000 \quad 0000$$

$$\mathbf{x}_3 = 111 \quad 0000 \quad 0000 \quad 1111 \quad 0000 \quad 0000$$

$$\mathbf{x}_4 = 111 \quad 0000 \quad 0000 \quad 0000 \quad 1111 \quad 0000$$

$$\mathbf{x}_5 = 111 \quad 0000 \quad 0000 \quad 0000 \quad 0000 \quad 1111$$

These give us 5 linearly independent elements which (apart from permutations of columns) must belong to any close-packed (23, 7) code. Now, as we are looking for a group code, we will try to find seven more linearly independent elements, guided by Theorems 2 and 3.

We have already written down the 5 points having 1 in the first 3 positions. Now, by Theorem 3, there are precisely 5 code points (of norm 7) having 1 in the 3 positions (say) 1, 2, 4. Of these we have so far only one, \mathbf{x}_1 , and so must be able to adjoin 4 more (if a (23, 7) code with 2^{11} points is to exist at all). In the construction we are guided by Theorem 2, which tells us that every set of four places must be filled with 1's exactly once. Hence, we add the expansions

$$\mathbf{x}_6 = 110 \quad 1000 \quad 1000 \quad 1000 \quad 1000 \quad 1000$$

$$\mathbf{x}_7 = 110 \quad 1000 \quad 0100 \quad 0100 \quad 0100 \quad 0100$$

$$\mathbf{x}_8 = 110 \quad 1000 \quad 0010 \quad 0010 \quad 0010 \quad 0010$$

$$\mathbf{x}_9 = 110 \quad 1000 \quad 0001 \quad 0001 \quad 0001 \quad 0001$$

Now, we could go on and adjoin new points, e.g., those having 1's in the positions 1, 2, 5 of which only one (\mathbf{x}_1) as at present in our collection. It is here, however, that some calculation is unavoidable if we seek to get linearly independent elements. Note that already $\mathbf{x}_1 + \mathbf{x}_3 + \mathbf{x}_4 + \mathbf{x}_5 = \mathbf{x}_6 + \mathbf{x}_7 + \mathbf{x}_8 + \mathbf{x}_9$ so that one of our 9 points must be discarded. In any case, however, it is clear that one by one, all the points of norm 7 can be written down until 11 linearly independent ones are obtained. Theorem 6 is a special case of the following:

• Theorem 7

Let $e \geq 3$ be odd, and let E denote the integer $e!(1 + 1/3 + \dots + 1/e)$. If a close-packed $(n, 2e+1)$ code exists, then one of the numbers $n+1$ or

$$e! \left\{ 1 + \frac{n(n-1)}{3!} + \dots + \frac{n(n-1) \dots (n-e+2)}{e!} \right\}$$

is a divisor of $2^a \cdot B$. Here 2^a denotes the largest power of 2 dividing E , and B denotes the largest odd divisor of $e!$.

Proof

If a close-packed code exists, $1 + \binom{n}{1} + \dots + \binom{n}{e}$ is a power

of 2, say 2^k . Now, when e is odd, this expression (considered as a polynomial in n) is factorable, since it vanishes for $n = -1$. It factors into

$$(1+n) \left\{ 1 + \frac{n(n-1)}{3!} + \dots + \frac{n(n-1) \dots (n-e+2)}{e!} \right\} = (1+n) \frac{P(n)}{e!},$$

say, and we have

$$(n+1)P(n) = e!2^k.$$

Now, $P(n) = P(-1) + (n+1)Q(n)$ where Q is some polynomial with rational integer coefficients. Hence any common divisor of $n+1$ and $P(n)$ divides $P(-1)$ which is seen to be E . Write $E = 2^a A$, where A is odd, and $e! = 2^b B$, where B is odd. Suppose $n+1$ fails to divide $2^a B$. Since any odd divisor of $n+1$ must divide B , this can only happen if 2^{a+1} divides $n+1$. In this case, 2^{a+1} cannot also divide $P(n)$, or else it would divide E , an impossibility. Therefore $P(n)$ divides $2^a B$, completing the proof.

Remark 1

It seems doubtful that $P(n)$ could actually divide $2^a B$ in practice, since it is probably greater than this number when $n \geq 2e+1$; however, to substantiate this argument one would need an estimate on the largest power of 2 dividing E . This seems in general to be a difficult number-theoretic problem, except in the case when $e \equiv 1 \pmod{4}$, where we have:

Corollary 1

If a close-packed $(n, 2e+1)$ code exists and e is of the form $4r+1$, then $n+1$ is a divisor of $e!$.

Proof of corollary: In this case,

$E = e!(1 + 1/3 + \dots + 1/e) = 2^b(B + B/3 + \dots + B/e) = 2^b B'$ where B' is odd, being the sum of an odd number of odd summands. Thus 2^b is the highest power of 2 dividing E , and Theorem 7 asserts that either $n+1$ or $P(n)$ divides $e!$. But the latter alternative cannot happen since $P(n) > P(0) = e!$ for all $n \geq 1$.

As illustrations, consider first $e=3$. Here $E=8$, and by Theorem 7 either $n+1$ or n^2-n+6 divides 24; this readily yields Theorem 6. Again, take $e=5$. By Corollary 1, $n+1$ divides 120, and we have:

Corollary 2

If a close-packed 5-error correcting code on n digits exists, then n has one of the values 11, 14, 19, 23, 29, 39, 59, or 119.

It might be of interest to further investigate these values

of n , and see whether $1 + \binom{n}{1} + \dots + \binom{n}{5}$ is actually a power

of 2 for any $n > 11$ among these numbers. If the answer is affirmative, one can then systematically search for a corresponding code along the lines indicated earlier. This program could easily be carried out by a digital computer.

Remark 2

Theorem 7 shows that if $e > 1$ is odd, there are at most finitely many close-packed e -error correcting codes. We also

know this to be the case for $e=2$, by Theorem 5. We now show that this is true for all $e > 1$, but unfortunately in a manner which is non-constructive, i.e., allows no estimate of the possible number of close-packed codes.

• Theorem 8

If $e \geq 2$, the number of close-packed $(n, 2e+1)$ codes is finite.

The proof is a simple consequence of a deep result of C. L. Siegel from the theory of numbers.

Lemma (Siegel)

Let $f(x)$ be any polynomial which takes integer values when x is an integer. Then, unless $f(x)$ is a constant times a power of a linear polynomial, the largest prime factor of $f(n)$ increases without limit as $n \rightarrow \infty$.

To deduce Theorem 8 from this we have simply to verify that $f(x)$, defined by

$$f(x) = 1 + \binom{x}{1} + \dots + \binom{x}{e}, \quad (30)$$

is not a power of a linear polynomial if $e \geq 2$; then $f(n)$ has a prime factor > 2 for n sufficiently large, and so cannot be a power of 2. But suppose $f(x) = a(b+cx)^e$, where a, b, c are rational (as we may obviously assume).

Then $1 = f(0) = ab^e$, so we may write $f(x) = (1+rx)^e$ where $r = c/b$ is rational. Setting $x=1$ we get $2 = (1+r)^e$, so that $\sqrt[e]{2}$ is rational. This contradiction establishes the theorem.

4. Lower bounds for $B(n, d)$

To motivate the considerations of this section, let us note the following:

If an (n, d) group code \mathcal{C} is maximal, then every point of \mathcal{G}_n has distance $\leq d-1$ from some code point.

Note that a maximal group code is not necessarily a maximal code, since $A(n, d)$ code points will in general not be achieved with a group code. Still, the above proposition asserts that it is relatively maximal, i.e., one cannot adjoin new points to \mathcal{C} and get a larger (n, d) code. The proof is immediate: if there were a point w in \mathcal{G}_n at distance $\geq d$ from all points of \mathcal{C} , the set \mathcal{C}' of all points $x, x+w$ where x ranges over \mathcal{C} would be an (n, d) group code having twice as many points as \mathcal{C} .

The preceding proposition can be viewed as the case $r=0$ of the following:

• Theorem 9

If \mathcal{C} is an (n, d) group code having 2^k points and there exists a point w in \mathcal{G}_n at distance $\geq d-r$ from all points of \mathcal{C} , then there exists an $(n+r, d)$ group code having 2^{k+1} points.

Proof

Let \mathcal{C}' be the set of all x in \mathcal{C} , each augmented by r 0's, plus the set of all $x+w$ (x in \mathcal{C}), each augmented by r 1's. This \mathcal{C}' has the required properties. We are mainly interested in the case $r=1$, for which the following theorem, in a certain sense converse to Theorem 9, holds:

• Theorem 10

Let \mathcal{C} be a maximal (n, d) group code and suppose $B(n-1, d) < B(n, d)$. Then there exist a maximal $(n-1, d)$ group code \mathcal{C}' and a point x in \mathcal{G}_{n-1} at distance $\geq d-1$ from all points of \mathcal{C}' , such that \mathcal{C} arises from \mathcal{C}' and x by the adjunction procedure of the previous theorem.

Proof

Let \mathcal{C}' be that subset of \mathcal{C} whose n th coordinate is 0. Then \mathcal{C}' is a subgroup of \mathcal{C} and the subset \mathcal{C}'' of \mathcal{G}_{n-1} obtained by dropping the last 0 from each point of \mathcal{C}' is an $(n-1, d)$ group code. Furthermore, it has $B(n, d)/2$ points since \mathcal{C}' has half as many points as \mathcal{C} . To see this, note first that \mathcal{C}' is not all of \mathcal{C} , or else by dropping the last digit we would obtain an $(n-1, d)$ code with $B(n, d)$ points, contrary to assumption. Hence, the complement \mathcal{C}'' of \mathcal{C}' in \mathcal{C} is not empty, and any pair of points of \mathcal{C}'' (having 1 in the last digit) differ by an element of \mathcal{C}' so that they all belong to the same coset of \mathcal{C}' in \mathcal{C} , i.e., \mathcal{C}' has index 2 in \mathcal{C} . Let now y be a point of \mathcal{G}_{n-1} gotten by choosing any point of \mathcal{C}'' and deleting the last digit. Then y has distance $\geq d-1$ from all of \mathcal{C}' and we may reconstruct \mathcal{C} from \mathcal{C}' and y by the procedure of the last theorem (for the case $r=1$).

From these two theorems we see that it is in principle possible to construct a maximal (n, d) code for each n successively, by repeated adjunction. $B(n, d)$ will double in passage from n to $n+1$ except for values of n where the maximal code we have constructed "saturates" the n -cube, i.e., where it is not merely maximal (no point at distance $\geq d$ from all code points) but it fills out the cube so densely that no point even has distance $\geq d-1$ from all code points. Viewed in this context, the highest possible saturation of the n -cube is achieved when a close-packed (n, d) code exists ($d=2e+1$) for then there is no point of \mathcal{G}_n at distance $\geq e+1$ from all code points.

As a first application of these ideas it is instructive to see what happens in Hamming's case where $d-1=e+1$ so that the two extremes of saturation coincide. By Theorem 9 (for $r=1$), $B(n+1, 3)=2B(n, 3)$ unless every maximal $(n, 3)$ group code has no point at distance 2 from it, i.e., unless it is close-packed (hence, has $2^n/(n+1)$ points). But this cannot happen unless n has the form 2^k-1 . We thus have the result of Hamming: $B(n, 3)$ continues to double as n increases, except when n passes through one of the values 3, 7, 15, 31..., when it stays the same. Actually we have shown somewhat more, namely that when n does not have one of these exceptional values, any maximal (n, d) group code can be extended to a code of order $n+1$ with twice as many points. Also, by Theorem 10, $B(n, 3)=B(n+1, 3)$ when n takes one of the exceptional values, or else there would exist a maximal $(n, 3)$ code having a point at distance 2 from all its points, violating the close-packed property.

Now, the general case ($e \geq 2$) is not so simple as this; for, whereas in the Hamming case one could proceed by adjunction in any manner whatever and at each step arrive at a maximal code, in the general case it is possible that two different adjunctions at some stage will lead to two codes, one of which saturates the cube and the other not, so a

choice must be made if one is to keep on obtaining maximal codes by this procedure. Further investigation of these matters would seem worthwhile.

We can, in any case, use this method to obtain a fairly good lower bound for $B(n, d)$.

Lemma

$$\text{If } B(n, d) < \frac{2^n}{1 + \binom{n}{1} + \dots + \binom{n}{d-2}}$$

then $B(n+1, d) = 2B(n, d)$.

Proof

Let \mathcal{C} be a maximal (n, d) group code, and associate to each x of \mathcal{C} the "sphere" \mathcal{S}_x consisting of all points of \mathcal{G}_n having distance $\leq d-2$ from x . The number of points in each \mathcal{S}_x is $1 + \binom{n}{1} + \dots + \binom{n}{d-2}$, and there are $B(n, d)$ spheres. If the hypothesis of the lemma holds there is at least one point in \mathcal{G}_n lying in none of these \mathcal{S}_x , i.e., having distance $\geq d-1$ from all points of \mathcal{C} . Hence, by Theorem 9, $B(n+1, d) = 2B(n, d)$.

• Theorem 11

For an infinite sequence of n

$$B(n, d) \geq \frac{2^n}{1 + \binom{n}{1} + \dots + \binom{n}{d-2}} \quad (31)$$

Proof

If (31) failed to hold from some n on, then by the lemma $B(n, d)$ would continue doubling from that value n_0 onward, giving $B(n, d) = 2^{n-n_0} B(n_0, d)$ and letting $n \rightarrow \infty$, the inequality (1) of Hamming is violated.

Remark

It is possible to replace the right side of (31) by a somewhat smaller expression such that the resulting inequality holds for all n , and also to give information about the density of n for which (31) holds.

Of course, from the basic property of a maximal code we have the weaker inequality

$$B(n, d) \geq \frac{2^n}{1 + \binom{n}{1} + \dots + \binom{n}{d-1}} \quad (32)$$

since, if (32) failed to hold, the spheres of radius $d-1$ about the points of a maximal group code would not exhaust \mathcal{G}_n , hence there would be a point of \mathcal{G}_n at distance $\geq d$ from all these code points, which by the introductory remarks of Section 4 would lead to a larger (n, d) group code.

The orders of magnitude of the Hamming upper bound, and the lower bounds (31) and (32) are, respectively,

$$\frac{2^n}{n^e}, \frac{2^n}{n^{2e-1}} \text{ and } \frac{2^n}{n^{2e}}.$$

Thus, even for $e=2$ there is a significant gap between Hamming's upper bound and the lower bound (31), and further results would be of interest.

5. The p -ary codes

Let p denote any prime number (the reasons why we require this condition will become clear shortly). Suppose we consider now the set $\mathfrak{S}^{(p)}_n$ of n -tuples of symbols each chosen from a set of p (which we may of course designate by $0, 1, \dots, p-1$, the residues of the prime p). $\mathfrak{S}^{(p)}_n$ has p^n "points" which no longer permit such a simple geometric interpretation when $p > 2$. We may define the distance between two points as the number of coordinates in which they are different, and by $A^{(p)}(n, d)$ the maximum number of points of $\mathfrak{S}^{(p)}_n$ which may have minimum mutual distance d . Then, by an argument similar to that of Section 1 we have the analog of Hamming's inequality

$$A^{(p)}(n, 2e+1) \leq \frac{p^n}{1 + (p-1)\binom{n}{1} + \dots + (p-1)\binom{n}{e}}. \quad (33)$$

We may define a *group* or systematic code as before to be an (n, d) code which is a group under digit-wise addition (mod p) and introduce the corresponding notation $B^{(p)}(n, d)$. The distance between the points x and y of a group code is $\|x - y\|$ where $\|w\|$ is defined to be the number of non-zero digits in w . One sees immediately that a subgroup of $\mathfrak{S}^{(p)}_n$ is an (n, d) code if and only if all its non-zero elements have at least d non-zero digits.

It is an interesting fact that virtually all of the preceding results can be carried over to these p -ary codes. We will content ourselves with only one theorem, which shows that the analogues of Hamming's single-error correcting codes exist.^{1, 5, 11}

• Theorem 12

$$B^{(p)}(n, 3) = \frac{p^n}{1 + (p-1)n} \quad (34)$$

whenever the right side is an integer, i.e., whenever n is one of the numbers

$$\frac{p^k - 1}{p - 1} \quad (k = 2, 3, \dots); \quad (35)$$

moreover, for all n , $B^{(p)}(n, 3)$ is equal to the greatest power of p not exceeding the right side and $B^{(p)}(n+1, 3) = pB^{(p)}(n, 3)$ except when n is one of the numbers (35), in which case $B^{(p)}(n+1, 3) = B^{(p)}(n, 3)$.

Proof

We can actually write down the codes in question. Suppose

$$n = \frac{p^k - 1}{p - 1};$$

we wish to construct an $(n, 3)$ group code with p^{n-k} points. Let us begin by writing out a basis for the group $\mathfrak{S}^{(p)}_{n-k}$:

$$\begin{array}{l} n-k \\ \underbrace{1 \ 0 \dots 0} \\ 0 \ 1 \dots 0 \\ 0 \ 0 \dots 1 \end{array} \quad (36)$$

This consists of just $n-k$ "unit vectors" which, when all

possible sums of them are taken, precisely generate all elements of $\mathfrak{S}^{(p)}_{n-k}$. Suppose now we can adjoin k additional digits to these such that the resulting points x_1, \dots, x_{n-k} have the following property: every linear combination

$$a_1 x_1 + \dots + a_{n-k} x_{n-k}, \quad (37)$$

where the a_i are integers (mod p), not all 0, has at least 3 non-zero digits. Then we will indeed have constructed a basis for the code with the required properties. Now, let us consider the class \mathfrak{K} of all k -tuples of integers (mod p) at least two of which are distinct from 0. Their number is precisely $p^k - 1 - k(p-1)$ since we are excluding precisely $1 + k(p-1)$ from the total of p^k . Let us call two of these k -tuples x, y *equivalent* if $ax + by = 0$ (the k -tuple $00\dots 0$) for some $a, b \not\equiv 0 \pmod{p}$. This is an equivalence relation: $x - x = 0$, symmetry is obvious, and $ax + by = 0, cy + dz = 0$, implies $(ea)x + dz = 0$ where e is so chosen that $eb \equiv -c \pmod{p}$. Such an e exists because the integers (mod p) form a field, and division by non-zero elements is possible. (Here is where the assumption that p is a prime enters essentially). Now, if x is in \mathfrak{K} , the equivalent elements in \mathfrak{K} are the set of solutions y of $ax + by = 0$ for $a, b \not\equiv 0$, i.e., of $y = b^{-1}ax$ (where $b^{-1} \equiv 1 \pmod{p}$), i.e., the equivalent elements are simply multiples of x by the numbers $1, 2, \dots, p-1$. Since there are $p-1$ elements in each equivalence class the number of classes is

$$\frac{p^k - 1 - k(p-1)}{p-1} = n - k.$$

Hence, we may adjoin mutually non-equivalent elements of \mathfrak{K} to the "unit vectors" in the array (36). This is a basis for the required code; for consider a linear combination (37). If only one a_i is distinct from zero (37) has the form $a_i x_i$ which has by construction at least three non-zero digits. If precisely two a_i are $\not\equiv 0$, we have $a_i x_i + a_j x_j$ which has precisely two non-zero digits among the first $n-k$, and at least one among the last k since we chose inequivalent k -tuples. Finally, if $l \geq 3$ of the a_i are $\not\equiv 0$, the first $n-k$ digits of (37) already contain l non-vanishing digits.

The other details of the proof can easily be deduced by the (appropriately generalized) line of reasoning of Section 4. In fact, the existence of the above code could also be proved by this reasoning precisely as we derived Hamming's results. For instance, Theorem 9 with $r=1$ (extended in the obvious way) tells us that

$$B^{(p)}(n, 3) = pB^{(p)}(n, 3) \text{ when } n \neq \frac{p^k - 1}{p - 1}, \quad (38)$$

and Theorem 10 says that (38) cannot hold if n has one of these excepted values, or else the code constructed above would have a point at distance 2 from all its elements, implying that inequality holds in (34), contradicting what we just proved.

The above proof is completely constructive: namely, in adjoining the k -tuples to the "unit vectors" we simply pick any k -tuple from the set \mathfrak{K} for our first choice; discard from \mathfrak{K} the $p-1$ non-zero multiples of this, and pick any of the remaining k -tuples for our second choice; discard the multiples of this one in turn, et cetera, until all $n-k$ have been chosen.

In view of the ease with which these codes may be constructed and their efficiency, it might be worth considering their use in connection with a binary channel, by encoding the p symbols into binary. If p is chosen less than but near 2^m for some integer m , the encoding can be done with m binary digits per p -ary symbol, and with little waste, e.g., $p=7$, $m=3$; $p=31$, $m=5$.

References

1. M. J. E. Golay, "Notes on Digital Coding," *Proc. IRE*, **37**, 657 (1949).
2. M. J. E. Golay, "Binary Coding" *IRE Transactions PGIT-4*, **23**, (1954).
3. R. W. Hamming, "Error Detecting and Error Correcting Codes," *Bell System Tech. J.*, **29**, 147 (1950).
4. E. Hecke, *Algebraische Zahlen*, Chelsea, N. Y.
5. C. Y. Lee, "Some Properties of Nonbinary Error-Correcting Codes," *IRE Transactions IT-4*, **77** (1958).
6. S. P. Lloyd, "Binary Block Coding," *Bell System Tech. J.*, **36**, 517 (1957).
7. R. E. A. C. Paley, *Jour. Math. and Phys.*, **12**, 311 (1933).
8. M. Plotkin, "Minimum Distance Codes," Moore School of Electrical Engineering Report. See also P. Elias, "Coding for Noisy Channels" (Appendix), *IRE National Convention Record*, Pt. 4, 37-46 (1955).
9. D. Slepian, "A Class of Binary Signalling Alphabets," *Bell System Tech. J.*, **35**, 203 (1956).
10. E. N. Gilbert, "A Comparison of Signalling Alphabets," *Bell System Tech. J.*, **31**, 504 (1952).
11. W. Ulrich, "Non-Binary Error Correction Codes," *Bell System Tech. J.*, **36**, 1341 (1957).

Received July 29, 1958