# 応用数学特論 II (集中講義)

## Day 4 Error-Correcting Codes

盧 暁南 (山梨大学)
Xiao-Nan LU (University of Yamanashi)

Aug. 30, 2021
Kobe University

# Outline

1. **Ulam's game: A toy example for error-correction**

2. Linear codes

3. Hamming codes and projective geometry

4. Perfect codes and MDS codes

# Ulam's game

## Ulam's game (Rényi, 1961; Ulam, 1976)

Given $N$, guess a number $n$ ($0 \leq n < N$) by asking a series of yes-or-no questions. The respondent is permitted to lie at most once.



Stanisław Ulam (1909–1984)



Adventures of a Mathematician

(1st ed. published in 1976)



Adventures of a Mathematician

(film released in 2020)

# Example of Ulam's game with $N = 16$

- Please think of a number $n$ with $0 \leq n \leq 15$.
- Please answer the following questions.
  - You are permitted to lie to at most one question.
  - Truth-telling for all the questions is allowed.

1. Is $n \geq 8$?
2. Is $n \in \{4, 5, 6, 7, 12, 13, 14, 15\}$?
3. Is $n \in \{2, 3, 6, 7, 10, 11, 14, 15\}$?
4. Is $n$ odd?
5. Is $n \in \{1, 2, 4, 7, 9, 10, 12, 15\}$?
6. Is $n \in \{1, 2, 5, 6, 8, 11, 12, 15\}$?
7. Is $n \in \{1, 3, 4, 6, 8, 10, 13, 15\}$?

| | | |
|---|---|---|
| 1 | Yes | 1 |
| 2 | Yes | 1 |
| 3 | No | 0 |
| 4 | Yes | 1 |
| 5 | No | 0 |
| 6 | Yes | 0 |
| 7 | No | 0 |

- According to your answers, I guess

$$n = 5$$

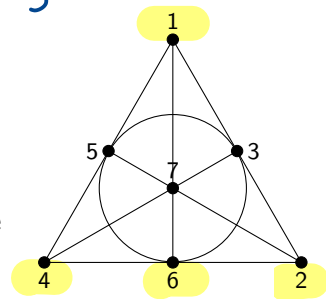I guess that you lied to question

No. 1.

# Catch the liar

$$y = (1 1 0 1 0 0 0)$$
$$x = (1 1 0 1 0 1 0)$$

- The answers can be represented by a binary vector $\mathbf{x}$ in $\{0,1\}^7$, where 1=yes, 0=no.

- The number of 1s is called the weight (重み) of $\mathbf{x}$, denoted by $\mathrm{wt}(\mathbf{x})$.

$$\mathrm{wt}(x) = 4. \quad \mathrm{wt}(y) = 3$$

1. If $\mathrm{wt}(\mathbf{x}) = 0$, no lie.

2. If $\mathrm{wt}(\mathbf{x}) = 1$, the lie = the position of 1.

3. If $\mathrm{wt}(\mathbf{x}) = 2$, two positions of 1s are lying in a unique line $\ell$ in Fano plane. The lie = the third point on $\ell$.

4. If $\mathrm{wt}(\mathbf{x}) = 3$, if three positions of 1s are lying in a unique line, then no lie. Otherwise, the four positions of 0s contains three point which form a line and one other point $P$. The lie = $P$.

5. If $\mathrm{wt}(\mathbf{x}) \geq 4$, apply the above rules for positions of 0s.

# Catch the liar

- The answers can be represented by a binary vector $\mathbf{x}$ in $\{0,1\}^7$, where 1=yes, 0=no.
- The number of 1s is called the weight (重み) of $\mathbf{x}$, denoted by $\mathrm{wt}(\mathbf{x})$.

- $\mathbf{x} = (0,0,0,0,0,0,1)$. The lie $= 7$.
- $\mathbf{x} = (0,1,0,0,0,0,1)$. There is a line $\{2,5,7\}$, so the lie $= 5$.
- $\mathbf{x} = (0,1,0,1,0,0,1)$. $\{2,4,7\}$ does not form a line. $\{1,3,5,6\}$ contains a line $\{3,5,6\}$. So the lie $= 1$.
- $\mathbf{x} = (0,1,0,1,0,1,1)$. $\{1,3,5\}$ does not form a line. $\{2,4,6,7\}$ contains a line $\{2,4,6\}$. So the lie $= 7$.
- $\mathbf{x} = (0,1,0,1,1,1,1)$. There is a line $\{1,2,3\}$, so the lie $= 2$.
- $\mathbf{x} = (1,1,0,1,1,1,1)$. The lie $= 3$.

## Correct the error

- Flip the bit ($0 \to 1$, $1 \to 0$) on the liar position. Denote the corrected vector by $\tilde{\mathbf{x}}$.
- Take the first four bits of $\tilde{\mathbf{x}}$ (a binary number) and convert it to decimal.

- $\mathbf{x} = (0, 0, 0, 0, 0, 0, 1)$. The lie = 7. $(0000)_2 \to (0)_{10}$.
- $\mathbf{x} = (0, 1, 0, 0, 0, 0, 1)$. The lie = 5. $(0100)_2 \to (4)_{10}$.
- $\mathbf{x} = (0, 1, 0, 1, 0, 0, 1)$. The lie = 1. $(1101)_2 \to (13)_{10}$. $= 8 + 4 + 1 = 13$
- $\mathbf{x} = (0, 1, 0, 1, 0, 1, 1)$. The lie = 7. $(0101)_2 \to (5)_{10}$.
- $\mathbf{x} = (0, 1, 0, 1, 1, 1, 1)$. The lie = 2. $(0001)_2 \to (1)_{10}$.
- $\mathbf{x} = (1, 1, 0, 1, 1, 1, 1)$. The lie = 3. $(1111)_2 \to (15)_{10}$.

$x = (1 1 0 1 0 1 0)$. $lie = 1$. $\tilde{x} = (0 1 0 1 0 1 0)$. $n = 4 + 1 = 5$.

# Aim of today's lecture

## The most important aim of today's lecture

To understand the "coding theory" for winning Ulam's game.

- Basic tools: linear algebra (matrix multiplication, linear mapping), finite field (mainly $\mathbb{F}_2$)
- Basic ideas: finite projective geometry, combinatorial designs

# Outline

# Outline

# Codes

- A code (符号) $\mathcal{C}$ of length $n$ over an alphabet $\Omega$ is a subset of vectors in

$$\Omega^n = \{(x_1, x_2, \ldots, x_n) : x_i \in \Omega, 1 \le i \le n\}$$

- The vectors in a code $\mathcal{C}$ are called codewords (符号語).
- Usually, we consider $\Omega$ to be a finite field $\mathbb{F}_q$ or finite ring (e.g. $\mathbb{Z}_4$ for DNA encoding).
- For most applications in electronic communications, $\Omega = \mathbb{F}_2 = \{0, 1\}$.

# Hamming distance

- As the "space" of codes, consider the notion of distance in $\Omega^n$.

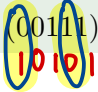### Hamming distance

The Hamming distance (ハミング距離) between any two codewords $\mathbf{x} = (x_1, \ldots, x_n)$, $\mathbf{y} = (y_1, \ldots, y_n)$ is the number of positions where they differ, i.e.

$$d_H(\mathbf{x}, \mathbf{y}) = \#\{i : x_i \neq y_i, 1 \leq i \leq n\}.$$

### Example (Hamming distance)

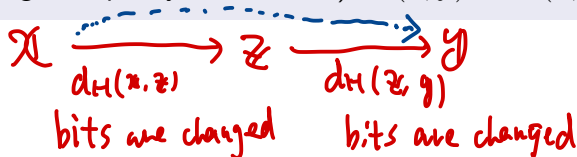For $\mathbf{x} = (00111)$, $\mathbf{y} = (10101)$, we have $d_H(\mathbf{x}, \mathbf{y}) = 2$.

# Hamming distance is a distance

> **Proposition** $d_H(x,y)$
>
> The Hamming distance is a metric (aka, distance function) on the space $\Omega^n$, namely, it satisfies the following conditions:
>
> 1. (nonnegativity; 非負性) $d_H(\mathbf{x}, \mathbf{y}) \geq 0$ for any $\mathbf{x}$ and $\mathbf{y}$.
> 2. (identity of indiscernibles; 同一律) $d_H(\mathbf{x}, \mathbf{y}) = 0$ if and only if $\mathbf{x} = \mathbf{y}$.
> 3. (symmetry; 対称性) $d_H(\mathbf{x}, \mathbf{y}) = d_H(\mathbf{y}, \mathbf{x})$ for any $\mathbf{x}$ and $\mathbf{y}$.
> 4. (triangle inequality; 三角不等式) $d_H(\mathbf{x}, \mathbf{y}) \leq d_H(\mathbf{x}, \mathbf{z}) + d_H(\mathbf{z}, \mathbf{y})$ for any $\mathbf{x}$, $\mathbf{y}$, and $\mathbf{z}$.

$x \xrightarrow{d_H(x,z)} z \xrightarrow{d_H(z,y)} y$

bits are changed    bits are changed

# Minimum distance

## Minimum distance

The minimum Hamming distance (最小ハミング距離) $d_{\mathcal{C}}$ of a code $\mathcal{C}$ is the smallest distance between any two distinct vectors,

$$d_{\mathcal{C}} = \min\{d_H(\mathbf{x}, \mathbf{y}) : \mathbf{x}, \mathbf{y} \in \mathcal{C}\}.$$

- If a code is of length $n$, has $M$ codewords, and minimum distance $d$, then the code is said to be an $(n, M, d)$ code.

# Linear codes over $\mathbb{F}_q$

- Now suppose $\Omega = \mathbb{F}_q$, the finite field of order $q$, where $q$ is a prime power.

### Linear code

$\mathcal{C} \subseteq \mathbb{F}_q^n$ is a linear code (線形符号) if the following conditions hold:

- If $\mathbf{v}$, $\mathbf{u} \in \mathcal{C}$ then $\mathbf{v} + \mathbf{u} \in \mathcal{C}$.
- If $\mathbf{v} \in \mathcal{C}$ and $\alpha \in \mathbb{F}_q$ then $\alpha \mathbf{v} \in \mathcal{C}$.

### Example (Repetition code)

Let $p$ be a prime. The following code

$$\mathcal{C} = \{(0, 0, \ldots, 0), (1, 1, \ldots, 1), \ldots, (p-1, p-1, \ldots, p-1)\}$$

is a linear code, called the repetition code (反復符号).

# Minimum weight of linear codes

## Hamming weight

The Hamming wight (ハミング重み) of a codeword $\mathbf{x} = (x_1, \ldots, x_n)$ in $\mathbb{F}_q^n$ is the number of its non-zero elements, i.e.

$$\mathrm{wt}(\mathbf{x}) = \#\{i : x_i \neq 0, 1 \leq i \leq n\}.$$

The minimum weight (最小重み) of a code is the smallest non-zero weight in the code.

zero codeword = $(0, 0, \ldots, 0)$

## Proposition

For a linear code $\mathcal{C}$, the minimum weight of the code is the minimum distance.

# Minimum weight of linear codes

## Hamming weight

The Hamming wight (ハミング重み) of a codeword $\mathbf{x} = (x_1, \ldots, x_n)$ in $\mathbb{F}_q^n$ is the number of its non-zero elements, i.e.

$$\text{wt}(\mathbf{x}) = \#\{i : x_i \neq 0, 1 \leq i \leq n\}.$$

The minimum weight (最小重み) of a code is the smallest non-zero weight in the code.

## Proposition

For a linear code $\mathcal{C}$, the minimum weight of the code is the minimum distance.

# Number of codewords in a linear code

## Proposition

Let $V$ be a vector space of dimension $k$ over $\mathbb{F}_q$, then $|V| = q^k$.

## Theorem

*A linear code of length $n$ over $\mathbb{F}_q$ is a subspace of $\mathbb{F}_q^n$. Hence, if $\mathcal{C}$ is a linear code over $\mathbb{F}_q$, then $|\mathcal{C}| = q^k$ for some $k$.*

$|\Omega| = 6$

- If a code of length $n$, with $M$ codewords, and minimum distance $d$ is an $(n, M, d)$ code.
- A linear code of length $n$, dim. $k$, and min. distance $d$ over $\mathbb{F}_q$ is an $[n, k, d]_q$ code.
- When $q = 2$, we usually omit the subscript $q$ and simple say an $[n, k, d]$ code.
- If the minimum distance $d$ is unknown, we say it is an $[n, k]_q$ code.

## Basis of a linear code

- $\mathcal{C}$ is an $[n,k]_q$ code $\iff$ $\mathcal{C}$ is a $k$-dimensional subspace of $\mathbb{F}_q^n$
- $\mathcal{C}$ has a basis (基底) of $k$ vectors. By elementary row / column operations (行列の行・列における基本変形), the matrix of basis can be transformed to $\begin{bmatrix} I_k & A \end{bmatrix}$.

### Example

Consider $\{(0,1,1,0),(1,1,1,1),(0,0,0,1)\}$ as a basis over $\mathbb{F}_2$. Then $(k=3)$

$$\begin{bmatrix} 0 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 \end{bmatrix} \xrightarrow[\text{(2)-(1)-(3)}]{\text{row}} \begin{bmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \xrightarrow[\substack{(2,\,1,\,4,\,3) \\ (1\ 2\ 3\ 4)}]{\text{column}} \left[\begin{array}{ccc:c} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{array}\right]$$

# Generator matrix of linear code

## Generator matrix

For an $[n,k]_q$ code $\mathcal{C}$, the matrix of a basis is a generator matrix (生成行列) of $\mathcal{C}$.

- A generator matrix of the form $\begin{bmatrix} I_k & A \end{bmatrix}$ is commonly used.

## Proposition

A linear code $\mathcal{C}$ with generator matrix $G$ can be obtained by

$$\mathcal{C} = \{\mathbf{x}G : \mathbf{x} \in \mathbb{F}_q^k\}.$$

Remark:

- In this lecture, we consider row vectors as codewords.
- If we consider column vectors as codewords, a generator matrix is of the form $\begin{bmatrix} I_k \\ A \end{bmatrix}$.

# Exercise 1: generator matrix of linear code

## Exercise 1

1. Find all the codewords of the linear code $\mathcal{C} \subseteq \mathbb{F}_2^4$ defined by the generator matrix

$$G = \begin{bmatrix} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$

2. Find another generator matrix of $\mathcal{C}$ that is different from $G$.

# Minimum distance decoding

ball: 球

sphere: 球面

## Hamming ball

An $n$-dimensional Hamming ball (ハミング球) of radius $r$ with center $\mathbf{v} \in \mathcal{C}$, where $\mathcal{C} \subseteq \mathbb{F}_q^n$ is a code, is the set of all the vectors in $\mathbb{F}_q^n$ having Hamming distance $\leq r$, i.e.

$$B_r(\mathbf{v}) = \{\mathbf{x} \in \mathbb{F}_q^n : d_H(\mathbf{x}, \mathbf{v}) \leq r\}$$

## Minimum distance decoding

When a vector $\mathbf{w}$ is received, it is decoded to the vector $\mathbf{v} \in \mathcal{C}$ that is closest to it. This method is called minimum distance decoding (最小距離復号) or nearest neighbor decoding.

# Error-correcting ability of linear codes

### Theorem

*Let $\mathcal{C}$ be an $[n, k, d]_q$ code, then $\mathcal{C}$ can detect $d - 1$ errors. If $d = 2e + 1$ then $\mathcal{C}$ can correct $e$ errors, and $\mathcal{C}$ is said to be an* e-error correcting code *(誤り訂正能力 $e$ の符号; $e$ 誤り訂正符号).*

$v_1, v \in \mathcal{C}.$

$d = 2e + 1$

$e$    $\geq d$    $e$

$v_1$    $v_2$

Hamming ball of radius $e$.

### Example (Repetition code)

The code $\mathcal{C}$ of length $n$ over $\mathbb{F}_q$ generated by $(1, 1, \ldots, 1)^\top$ has minimum distance of $n$. Hence, it can detect $n - 1$ errors and correct $\lfloor (n - 1)/2 \rfloor$ errors.

# Outline

# Dual code of linear code

- Let $\langle \mathbf{w}, \mathbf{v} \rangle$ denote the inner product (内積) of $\mathbf{w}, \mathbf{v} \in \mathbb{F}_q^n$.

## Dual code

For a linear code $\mathcal{C}$ over $\mathbb{F}_q$, let

$$\mathcal{C}^\perp = \{\mathbf{w} : \langle \mathbf{w}, \mathbf{v} \rangle = 0 \text{ for all } \mathbf{v} \in \mathcal{C}\}.$$

Then $\mathcal{C}^\perp$ is called the dual code (双対符号) of $\mathcal{C}$.

## Proposition

If $\mathcal{C}$ is an $[n, k]_q$ code, then $\mathcal{C}^\perp$ is an $[n, n-k]_q$ code.

## Example (Dual code of binary repetition code)

Let $\mathcal{C}$ be the code of length $n$ over $\mathbb{F}_2$ generated by $(1, 1, \ldots, 1)^\top$. Then $\mathcal{C}^\perp$ has dimension $n-1$ and consists of all vectors of length $n$ with even weight.

# Generator matrix of dual code

## Theorem

*If a linear code $\mathcal{C}$ is generated by $\begin{bmatrix} I_k & A \end{bmatrix}$, then $\mathcal{C}^{\perp}$ is generated by $\begin{bmatrix} -A^{\top} & I_{n-k} \end{bmatrix}$.*

## Exercise 1

❸ For the linear code $\mathcal{C} \subseteq \mathbb{F}_2^4$ generated by $G = \begin{bmatrix} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}$, produce all the codewords

of its dual code $\mathcal{C}^{\perp}$.

❹ Produce a generator matrix of $\mathcal{C}^{\perp}$. (Hint: you can use the above theorem. )

# Self-orthogonal and self-dual codes

### Self-orthogonal and self-dual codes

A code $\mathcal{C}$ is self-orthogonal (自己直交) if $\mathcal{C} \subseteq \mathcal{C}^{\perp}$ and it is self-dual (自己双対) if $\mathcal{C} = \mathcal{C}^{\perp}$.

### Exercise 1

⑤ Check whether $\mathcal{C}^{\perp}$ is a self-orthogonal code and briefly state the reason why it is or it is not.

# Parity check matrix and syndrome

- Let $\mathcal{C}$ be an $[n,k]_q$ code and suppose $\mathcal{C}^\perp$ has a generator matrix $H = \begin{bmatrix} I_{n-k} & B \end{bmatrix}$.
- Then $\mathbf{x} \in \mathcal{C} \iff H\mathbf{x}^\top = \mathbf{0}$. $\iff$ *$\mathcal{X}$ is orthogonal with codewords in $\mathcal{C}^\perp$.*
- This matrix $H$ is called the parity check matrix (パリティ検査行列) of $\mathcal{C}$.

---

**Syndrome**

*奇偶性*

Let $\mathcal{C}$ be a code in $\mathbb{F}_q$ with parity check matrix $H$. Then the syndrome (シンドローム [1].) of a vector $\mathbf{v} \in \mathbb{F}_q^n$ is $S(\mathbf{v}) = H\mathbf{v}^\top$. ← *column vector*

---

- For $\mathbf{x} \in \mathcal{C}$, we have $S(\mathbf{x} + \mathbf{e}) = H\mathbf{x}^\top + H\mathbf{e}^\top = H\mathbf{e}^\top$.

*受信列 v, 誤り e*

---

[1]In Chinese, "校正子" or "校験子"

# Binary Hamming code: The winning strategy for Ulam's game (1/4)

- Consider the linear code $\mathcal{C}$ generate by

$$G = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 & 1 \end{bmatrix}.$$

- Its dual code $\mathcal{C}^{\perp}$ has generator matrix

$$H = \begin{bmatrix} 0 & 1 & 1 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 & 1 \end{bmatrix},$$

which is the parity check matrix of $\mathcal{C}$.

# Binary Hamming code: The winning strategy for Ulam's game (2/4)

- Suppose that $n = 14$ is chosen, whose binary expression is $(1110)$.

- The codeword is

$$|\mathcal{C}| = 2^4 = 16$$

$$\mathbf{x} = \begin{bmatrix} 1 & 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & 0 & 0 & 0 & 0 \end{bmatrix}.$$

- Now sent the message $(1110000)$. If it is received as $\mathbf{x} = (1110000)$ then it is considered to be correct (with no lie). Because, the syndrome is $\mathbf{0}$, i.e.,

$$H\mathbf{x}^\top = \begin{bmatrix} 0 & 1 & 1 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 0 & 0 & 0 & 0 \end{bmatrix}^\top = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}.$$

# Binary Hamming code: The winning strategy for Ulam's game (3/4)

- Precisely, $\mathcal{C}$ has 16 codewords, in which $(1110000)$ is a codeword.

| | |
|---|---|
| $(0,0,0,0,0,0,0)$ | $(1,0,0,0,0,1,1)$ |
| $(0,0,0,1,1,1,1)$ | $(1,0,0,1,1,0,0)$ |
| $(0,0,1,0,1,1,0)$ | $(1,0,1,0,1,0,1)$ |
| $(0,0,1,1,0,0,1)$ | $(1,0,1,1,0,1,0)$ |
| $(0,1,0,0,1,0,1)$ | $(1,1,0,0,1,1,0)$ |
| $(0,1,0,1,0,1,0)$ | $(1,1,0,1,0,0,1)$ |
| $(0,1,1,0,0,1,1)$ | $(1,1,1,0,0,0,0)$ |
| $(0,1,1,1,1,0,0)$ | $(1,1,1,1,1,1,1)$ |

- In general, it is not smart to search through an entire code, because usually the codes contain an extremely large amount of codewords.

# Binary Hamming code: The winning strategy for Ulam's game (4/4)

- Assume that it is received as $\mathbf{x}' = (1110010)$. We have

$$H\mathbf{x}'^\top = \begin{bmatrix} 0 & 1 & 1 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 0 & 0 & 1 & 0 \end{bmatrix}^\top = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}.$$

  The syndrome is nothing but the 6th column of $H$. So we conclude that the 6th bit has an error (tells a lie).

- Assume that it is received as $\mathbf{x}' = (1010000)$. We have

$$H\mathbf{x}'^\top = \begin{bmatrix} 0 & 1 & 1 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 1 & 0 & 0 & 0 & 0 \end{bmatrix}^\top = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}.$$

  So we conclude that the 2nd bit has an error (tells a lie).

# Errors that $[7,4]_2$ Hamming code cannot correct

- Assume that it is received as $\mathbf{x}' = (1100010)$. We have

$$H\mathbf{x}'^{\top} = \begin{bmatrix} 0 & 1 & 1 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}^{\top} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}.$$

  So we found that the 5th bit has an error (tells a lie). Then we "correct" it to $(1100110)$. Wrong answer!

- If you told two lies or more, I could not win (by using the $[7,4]_2$ code).

# An illustration for a Hamming code



Figure from

Hamady, Micah, et al. Error-correcting barcoded primers allow hundreds of samples to be pyrosequenced in multiplex.

Nature Methods 5(3): 235–237., 2008.

# Outline

## Fano plane (revisit)

- Fano plane is $PG(2, \mathbb{F}_2)$ with 7 points $\{1, 2, \ldots, 7\}$ and 7 lines.
- The characteristic vectors (特性ベクトル) of lines are the red ones among the codewords of the $[7, 4]_2$ Hamming code.

$(0, 0, 0, 0, 0, 0, 0)$       $(1, 0, 0, 0, 0, 1, 1)$

$(0, 0, 0, 1, 1, 1, 1)$       $(1, 0, 0, 1, 1, 0, 0)$

$(0, 0, 1, 0, 1, 1, 0)$       $(1, 0, 1, 0, 1, 0, 1)$

$(0, 0, 1, 1, 0, 0, 1)$       $(1, 0, 1, 1, 0, 1, 0)$

$(0, 1, 0, 0, 1, 0, 1)$       $(1, 1, 0, 0, 1, 1, 0)$

$(0, 1, 0, 1, 0, 1, 0)$       $(1, 1, 0, 1, 0, 0, 1)$

$(0, 1, 1, 0, 0, 1, 1)$       $(1, 1, 1, 0, 0, 0, 0)$

$(0, 1, 1, 1, 1, 0, 0)$       $(1, 1, 1, 1, 1, 1, 1)$

## Some more geometric terms

- A $k$-arc (弧) is a set of $k$ points in a plane such that no three are collinear. An arc of maximal size is said to be an oval (卵形).
- The maximum size of an arc in a plane of order $n$ is $n+2$ (called a hyperoval) if $n$ is even and $n+1$ (called a oval) if $n$ is odd.
- Any four points in Fano plane that does not contain a line form a hyperoval.
- The complement of any line in Fano plane is a hyperoval.

For the $[7, 4]_2$ Hamming code $\mathcal{C}$,

$$\mathcal{C} = \{\text{lines in } \mathrm{PG}(2, \mathbb{F}_2)\} \cup \{\text{hyperovals in } \mathrm{PG}(2, \mathbb{F}_2)\}$$
$$\cup \{\text{all-one codeword}, \text{all-zero codeword}\}$$

# Why does the geometric winning strategy work?

1. If $\mathrm{wt}(\mathbf{x}) = 1$, the lie = the position of 1. [nearest neighbor to all-zero codeword]

2. If $\mathrm{wt}(\mathbf{x}) = 2$, two positions of 1s are lying in a unique line $\ell$ in Fano plane. The lie = the third point on $\ell$. [nearest neighbor to a line]

3. If $\mathrm{wt}(\mathbf{x}) = 3$, if three positions of 1s are lying in a unique line, then no lie. Otherwise, the four positions of 0s contains three point which form a line and one other point $P$. The lie = $P$. [a line or nearest neighbor to a hyperoval]

4. If $\mathrm{wt}(\mathbf{x}) = 4$, [a hyperoval or nearest neighbor to a line]

5. If $\mathrm{wt}(\mathbf{x}) = 5$, [nearest neighbor to a hyperoval]

6. If $\mathrm{wt}(\mathbf{x}) = 6$, [nearest neighbor to all-one codeword]

We were looking for the ~~nearest neighbor~~ to some codeword!

as nearest neighbor of $\mathbf{x}$

# Some more algebra ... Coset decomposition

## Proposition

For an $[n, k]_q$ code $\mathcal{C}$ and a vector $\mathbf{w} \in \mathbb{F}_q^n$,

$$\mathcal{C} + \mathbf{w} = \{\mathbf{x} + \mathbf{w} : \mathbf{x} \in \mathcal{C}\}.$$

is called a coset (コセット, 剩余類 [2]) of $\mathcal{C}$.

- For $\mathbf{w} \in \mathbb{F}_q^n$, $|\mathcal{C}| = |\mathcal{C} + \mathbf{w}|$.
- For $\mathbf{w}, \mathbf{u} \in \mathbb{F}_q^n$, either $\mathcal{C} + \mathbf{w} = \mathcal{C} + \mathbf{u}$ or $(\mathcal{C} + \mathbf{w}) \cap (\mathcal{C} + \mathbf{u}) = \emptyset$.
- There exist distinct vectors $\mathbf{w}_0, \ldots, \mathbf{w}_{q^{n-k}-1}$ such that

$$\mathbb{F}_q^n = \bigcup_{i=0}^{q^{n-k}-1} (\mathcal{C} + \mathbf{w}_i)$$

---

[2]In Chinese, "陪集".

# Final explanation for winning strategy of Ulam's game

## Proposition

Let $\mathcal{C}$ be an $[n, k]_q$ code with parity check matrix $H$. Let $\mathbf{v}, \mathbf{w} \in \mathbb{F}_q^n$. Then $S(\mathbf{v}) = S(\mathbf{w})$ if and only if $\mathbf{v}$ and $\mathbf{w}$ are in the same coset of $\mathcal{C}$.

- Each coset of the $[7, 4]_2$ Hamming code $\mathcal{C}$ has 16 codewords.
- Let $\mathbf{e}_i$ $(1 \leq i \leq 7)$ denote the vector in $\mathbb{F}_2^7$ whose $i$th entry is 1 and the others are 0.
- The cosets $\mathcal{C} + \mathbf{e}_i$ are distinct for distinct $i$. Hence,

$$\mathbb{F}_2^7 = \bigcup_{i=1}^{7} (\mathcal{C} + \mathbf{e}_i)$$

$e_1 = (1, 0, 0, \ldots, 0)$
$e_2 = (0, 1, 0, \ldots, 0)$

- For $\mathbf{x}' = \mathbf{x} + \mathbf{e}_i$ (one lie is told for the $i$th question),

$$H\mathbf{x}'^\top = H(\mathbf{x} + \mathbf{e}_i)^\top = H\mathbf{e}_i^\top = i\text{th column vector of } H.$$

# Summary on Hamming codes

## Binary Hamming code

$[7,4]_2$ code ⇐ (m=3 のとき)

A binary Hamming code is a $[n = 2^m - 1, k = 2^m - m - 1]_2$ code with parity check matrix $H$ whose $i$th ($1 \leq i \leq 2^m - 1$) column is the binary representation of $i$.

## Proposition

単一誤り訂正符号

A binary Hamming code corrects all single errors, i.e., its minimum distance is $d = 3$.

## Theorem (generalized Hamming code)

*Let $n = (q^m - 1)/(q - 1)$. Let $\mathbf{v}_1, \mathbf{v}_2, \ldots, \mathbf{v}_n$ be distinct points in $\mathrm{PG}(m - 1, \mathbb{F}_q)$, i.e. the vectors in $\mathbb{F}_q^m$. Let $H$ be the $m \times n$ matrix whose column vectors are $\mathbf{v}_1, \mathbf{v}_1, \ldots, \mathbf{v}_n$. Let $\mathcal{H}(m, q)$ be the code having $H$ as its parity check matrix. Then $\mathcal{H}(m, q)$ is called a generalized Hamming code and it is a $[\frac{q^m-1}{q-1}, \frac{q^m-1}{q-1} - m, 3]_q$ code.*

q-ary (q元 H.C.)

# Outline

# Sphere packing bound

## Proposition

- For any vector $\mathbf{v} \in \mathbb{F}_q^n$ there are $\binom{n}{s}(q-1)^s$ vectors in $\mathbb{F}_q^n$ that have Hamming distance $s$ from $\mathbf{v}$.

- For any vector $\mathbf{v} \in \mathbb{F}_q^n$ there are $\sum_{s=0}^{t} \binom{n}{s}(q-1)^s$ vectors in the sphere of radius $t$ centered at $\mathbf{v}$.

  *(Hamming ball)*

## Theorem (sphere packing bound)   球充填の限界式

*Let $\mathcal{C} \subseteq \mathbb{F}_q^n$ be a code with minimum weight $2t + 1$. Then*

$$|\mathcal{C}| \leq \frac{q^n}{\sum_{s=0}^{t} \binom{n}{s}(q-1)^s}.$$

$\longleftarrow |\mathbb{F}_q^n|$

$\longleftarrow$ # vectors in a Hamming ball

# Perfect codes

## Perfect codes

A code with equality in the sphere packing bound is said to be a perfect code.

- The generalized Hamming codes are perfect codes.



Figure from `https://en.wikipedia.org/wiki/Sphere_packing`

## Golay codes

- Golay's $[11, 6, 5]_3$ codes and Golay's $[23, 12, 7]_2$ codes are perfect codes (Golay, 1949).

- The binary Golay code is closely related to the Witt $4$- and $5$-designs and the Mathieu groups.



Reprinted from *Proc. IRE*, vol. 37, p. 657, June 1949.

M. J. E. Golay, Notes on digital coding. Proc. IEEE 37, p. 657, 1949.

by E. R. Berlekamp, the "best single published page" in coding theory.

# Singleton bound and MDS codes

## Theorem (Singleton bound)

*Let $\mathcal{C}$ be a code of length $n$ over an alphabet $\Omega$ of size $q$ with minimum Hamming distance $d$ and $q^k$ elements. Then*

$$d \le n - k + 1.$$

## MDS codes

A code with equality in the singleton bound is said to be a Maximum Distance Separable (MDS) code.

## MDS conjecture

If $\mathcal{C}$ is an $[n, k, n-k+1]_q$ MDS code then $n \le q + 1$.

# MDS codes and combinatorial structures

### Theorem

*A set of $s$ MOLS of order $q$ is equivalent to an $[s+2, 2, s+1]_q$ MDS code.*

### Theorem

*An $[n, k, n-k+1]_q$ MDS code is equivalent to an $n$-arc in $\mathrm{PG}(n-k-1, \mathbb{F}_q)$.*

# Homework assignments (レポート課題) for 4th day

## Exercise 1

▸ (1) (2) , ▸ (3) (4) , ▸ (5)

## Exercise 2

$$2^4 - 1 = 15 \qquad (m = 4)$$

1. Produce the parity check matrix for the Hamming code of length $15$.
2. Use the parity check matrix produced above to decode the vectors $(111001101101101)$ and $(001100110011010)$.

- You are encouraged to use computer programs.
- Deadline: 6th Sept., 23:59:59