

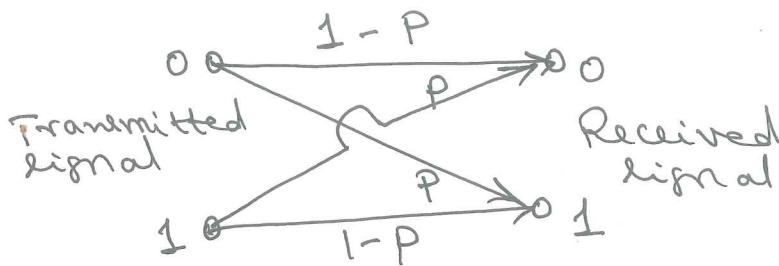
## Chapter - 10 : Coding Theory.

Introduction :- The algebraic coding theory was inspired by the fundamental paper of Claude Shannon (1948) along with results by Marcel Golay (1949) and Richard Hamming (1950). Since that time it has become an area of great interest where algebraic structures, probability, and combinatorics all play a role.

In this chapter, our coverage will be held to an introductory level as we seek to model the transmission of information represented by strings of the signals 0 and 1.

In digital communications, when information is transmitted in the form of strings of 0's and 1's, certain problems arise. As a result of "noise" in the channel, when a certain signal is transmitted a different signal may be received, thus calling the receiver to make wrong decision. Hence, we want to develop techniques to help us detect and perhaps even correct transmission errors. Only we can improve the chances of correct transmission there are no guarantees.

Our model uses a binary symmetric channel, as shown below.



The adjective binary appears because an individual signal is represented by one of the bits 0 or 1. When a transmitter sends the signal 0 or 1 in such a channel, associated with either signal is a probability  $p$  for incorrect transmission. When that probability  $p$  is same for both signals, the channel is called symmetric. Here we have probability  $p$  of sending 0 and having 1 received. The probability of sending signal 0 and having it received correctly is then  $1-p$ .

Ex: Consider a string  $c = 10110$ .

$c \in \mathbb{Z}_2^5$  group, formed from the direct product of five copies of  $(\mathbb{Z}_2, +)$ .

When sending each bit of  $c$  through the binary symmetric channel, we assume that the probability of incorrect transmission is  $p = 0.05$ , so that the probability of transmitting  $c$  with no errors is  $(0.95)^5 = 0.77$ .

NOTE: The transmission of each signal does not depend in any way on the transmissions of prior signals. Consequently, the probability of the occurrence of all of these independent events is given by the product of their individual probabilities.

What is the probability that the party receiving the five-bit message receives the string  $r = 00110$ ?

$\Rightarrow$  Here the original message is  $c = 10110$ , therefore the error occurs at first position.

The probability is 0.05 for the first bit and 0.95 for the remaining 4 bits. Therefore,  $(0.05)(0.95)^4$

$= 0.041$  is the probability of sending  $c = 10110$  and receiving  $r = 00110$ . By considering  $e = 10000$ , we can write  $c+e=r$ . That is,  $r$  is the sum of the original message and error pattern  $e=10000$ . Here,  $c+r=e$  or  $r+e=c$ , since,  $e, c, r \in \mathbb{Z}_2^5$  and  $1+1=0$  in  $\mathbb{Z}_2$ .

In transmitting  $c = 10110$ , the probability of receiving  $r = 00100$  is  $(0.05)(0.95)^2(0.05)(0.95) = 0.002$  hence, the occurrence of this multiple error is less likely to happen. (12)

What is the probability that  $r$  differs from  $c$  in exactly two places?

$$\binom{5}{2} (0.05)^2 (0.95)^3 = 0.021.$$

Two 1's & 3 zeros

From this we can conclude the following.

"Let  $c \in \mathbb{Z}_2^n$ , for the transmission of  $c$  through a binary symmetric channel with probability  $p$  of incorrect transmissions, (a) the probability of receiving  $r = c + e$ , where  $e$  is a particular error pattern consisting of  $k$  1's and  $(n-k)$  zero is  $\boxed{p^k (1-p)^{n-k}}$  (b) the probability that  $k$

errors are made in the transmission is  

$$\binom{n}{k} p^k (1-p)^{n-k}$$

The probability of making, at most one error in the transmission of  $c = 10110$  is  $(0.95)^5 + \binom{5}{1} (0.05)^1 (0.95)^4 = 0.977$ . Thus the chance of multiple errors in transmission will be considered negligible. This assumption is valid when  $p$  is very small. In a binary symmetric channel is considered "good" when  $p < 10^{-5}$ , or  $p < 1/2$ .

To improve the accuracy of transmission in a binary symmetric channel, certain types of coding schemes can be used where extra signals are provided.

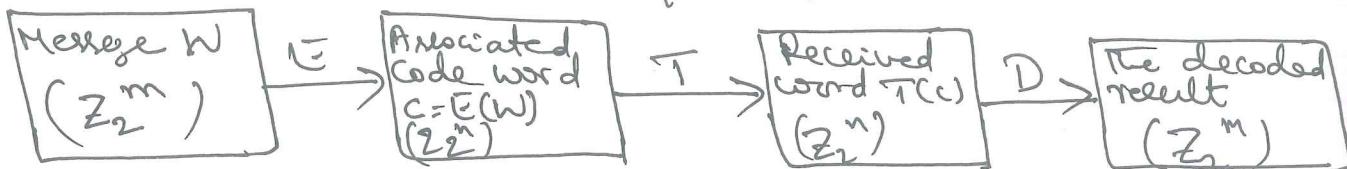
For  $m, n \in \mathbb{Z}^+$ , let  $n > m$ . Consider  $\phi \neq w \in \mathbb{Z}_2^m$ .

The set  $W$  consists of the messages to be transmitted. For each  $w \in W$ , append  $n-m$  extra bits to form the codeword  $c$ , where  $c \in \mathbb{Z}_2^n$ . This process is called encoding. This is represented by a function  $E: W \rightarrow \mathbb{Z}_2^n$ . Therefore  $E(w) = c$  and  $E(W) = C \subseteq \mathbb{Z}_2^n$ .

Since the function  $E$  simply appends extra bits to the distinct messages, the encoding is one-to-one.

After transmission of  $c$ , it is received as  $T(c)$ , where  $T(c) \in \mathbb{Z}_2^n$ . Unfortunately,  $T$  is not ~~one-to-one~~<sup>a function</sup> because  $T(c)$  may be different at different transmission times (noise changes with time).

Upon receiving  $T(c)$ , we want to apply a decoding function  $D: \mathbb{Z}_2^n \rightarrow \mathbb{Z}_2^m$  to remove the extra signals and to obtain the original message  $w$ . Ideally  $D \circ T \circ E$  should be the identity function on  $W$ , with  $D: C \rightarrow W$ . Since this cannot be expected, we seek functions  $E$  &  $D$  such that there is a high probability of correctly decoding the received word  $T(c)$  and recapturing the original message  $w$ .



The functions  $E$  &  $D$  are called encoding & decoding functions of an  $(n, m)$  block code.

Ex: Consider the  $(m+1, m)$  block code for  $m=8$ . (13)  
 Let  $W = \mathbb{Z}_2^8$ . For each  $w = w_1, w_2, \dots, w_8 \in W$ , define  
 $\mathbb{Z}_2^8 \rightarrow \mathbb{Z}_2^9$  by  $E(w) = w_1, w_2, \dots, w_8, w_9$ , where  $w_9 = \sum_{i=1}^8 w_i$   
 where the addition is modulo 2. For example,  
 $E(11001101) = 110011011$ , &  $E(00110011) = 001100110$   
 For all  $w \in \mathbb{Z}_2^8$ ,  $E(w)$  contains an even number of 1's. So for  $w = 11010110$  and  $E(w) = 110101101$ , if we receive  $T(c) = T(E(w))$  as 100101101, this has odd no. of 1's, hence we know that mistake happens in transmission. Hence, we are able to detect single errors in transmission. But, we don't know how to correct such errors.

The probability of sending the code word 110101101 and allowing at most one error in transmission is

$$(1-p)^9 + \binom{9}{1} p (1-p)^8 \quad \text{for } p=0.001, \text{ this gives}$$

$\underbrace{(1-p)^9}_{\substack{\text{All 9 bits are transmitted correctly}}} + \underbrace{\binom{9}{1} p (1-p)^8}_{\substack{\text{One bit changed in transmission and it is detected.}}}$ 
 $(0.999)^9 + \binom{9}{1} (0.001) (0.999)^8 = 0.999964$

If we detect an error and we are able to replay a signal back to the transmitter to repeat the transmission of the code word, and continue this process until the received word has an even number of 1's then the probability of sending the code word 110101101 and receiving the correct transmission is approximately 0.999964.

This  $(m+1, m)$  block code is called parity-check code. What is the drawback of this code?

E2: The  $(3m, m)$  block code [Triple Repetition Code] is one in which we can detect and correct single errors in transmission.

With  $m=8$  and  $W = \mathbb{Z}_2^8$ , we define

$$E: \mathbb{Z}_2^8 \rightarrow \mathbb{Z}_2^{24} \text{ by } E(w_1 w_2 \dots w_8) = w_1 w_1 \dots w_8 w_1 w_2 \dots w_8$$

The decoding function  $D: \mathbb{Z}_2^{24} \rightarrow \mathbb{Z}_2^8$  is carried out by the majority rule.

for example, if the  $w = 10110111$  then

$$c = E(w) = 10110111; 10110111; 10110111$$

If the received word  $T(c) = 10100011; 00110111; 10110110$ . In this we have 3 errors at position 4, 9 & 24.

Decode this  $T(c)$  by examining 1<sup>st</sup>, 9<sup>th</sup>, 17<sup>th</sup> bit and where which signal appears more times. Here it is 1, so the first bit is 1 in the decoded message. Similarly counting in this fashion for 2<sup>nd</sup>, 10<sup>th</sup> & 18<sup>th</sup> position, here all 3 are 0s hence the 2<sup>nd</sup> bit in decoded word is 0. By continue in this fashion we decoded  $T(c)$  back to  $w = 10110111$ .

What is the Draw back in this?

Even though, many errors are detected and corrected, but if many or 2-errors occur at the same position it can't be corrected / or incorrect transmission occurs.

What is the probability?  $p = 0.001$ .

$$(0.999)^3 + \binom{3}{1} (0.001) (0.999)^2 = 0.999997$$

$$\therefore (0.999997) = 0.999976$$

## Hamming Metric:

Consider a code  $C \subseteq \mathbb{Z}_2^4$ , where  $c_1 = 0111$  &  $c_2 = 1111 \in C$ . If both sender and receiver know the elements of  $C$ . If the transmitter sends  $c_1$  and the person receiving the code word receives  $T(c_1)$  as  $1111$ , then he or she feels that  $c_2$  was transmitted and makes the appropriate decision (wrong)  $c_2$  implied. Here, the two code words they are almost similar, they are close to each other, for they differ in only one component.

Definition: for each element  $x = x_1 x_2 \dots x_n \in \mathbb{Z}_2^n$ , where  $n \in \mathbb{Z}^+$ , the weight of  $x$ , denoted as  $\text{wt}(x)$ , is the number of components  $x_i$  of  $x$ , for  $1 \leq i \leq n$ , where  $x_i = 1$ . If  $y \in \mathbb{Z}_2^n$ , the distance between  $x$  &  $y$ , denoted by  $d(x, y)$ , is the number of components where  $x_i \neq y_i$ , for  $1 \leq i \leq n$ .

Ex: For  $n=5$ , let  $x=01001$  and  $y=11101$ ,

Then  $\text{wt}(x)=2$ ,  $\text{wt}(y)=4$  and  $d(x, y)=2$ .

$x+y = 10100$ , so  $\text{wt}(x+y)=2$ , In this example by chance  $d(x, y) = \text{wt}(x+y)$ .

for each  $1 \leq i \leq 5$ ,  $x_i + y_i$  contributes a count of 1 to  $\text{wt}(x+y) \Leftrightarrow x_i \neq y_i$ , whereas  $x_i, y_i$  contribute a count of 1 to  $d(x, y)$ .

$\therefore$  This is actually true for all  $n \in \mathbb{Z}^+$ , so  $\text{wt}(x+y) = d(x, y)$

Lemma: For all  $x, y \in \mathbb{Z}_2^n$ ,  $\boxed{\text{wt}(x+y) \leq \text{wt}(x) + \text{wt}(y)}$

$$x = x_1 \dots x_n$$

$$y = y_1 \dots y_n$$

$$x+y = z_1 z_2 \dots z_n$$

The distance function  $d$ , defined on  $\mathbb{Z}_2^n \times \mathbb{Z}_2^n$  satisfies the following for all  $x, y, z \in \mathbb{Z}_2^n$

- a]  $d(x, y) \geq 0$
- b]  $d(x, y) = 0 \iff x = y$
- c]  $d(x, y) = d(y, x)$
- d]  $d(x, z) \leq d(x, y) + d(y, z)$

Any function which satisfies all the above properties is called a distance function or metric, and we call  $(\mathbb{Z}_2^n, d)$  a metric space. Here  $d$  is often referred as Hamming metric:

Definition: (Sphere): for  $n, k \in \mathbb{Z}^+$  and  $x \in \mathbb{Z}_2^n$ , the sphere of radius  $k$  centered at  $x$  is defined as

$$S(x, k) = \{y \in \mathbb{Z}_2^n \mid d(x, y) \leq k\}$$

Q1: for  $n=3$  and  $x=110 \in \mathbb{Z}_2^3$

$$\therefore S(x, 1) = \{110, 010, 100, 111\}$$

$$S(x, 2) = \{110, 010, 100, 111, 000, 101, 011\}$$

What is the error detecting and correcting capabilities of Hamming metric?

Let  $E: W \rightarrow C$  be a encoding function with the set of messages  $W \subseteq \mathbb{Z}_2^m$  and the set of code words  $E(W) = C \subseteq \mathbb{Z}_2^n$ , where  $m < n$ . For  $k \in \mathbb{Z}^+$ , we can detect transmission errors of weight  $\leq k$  iff the minimum distance between code words is at least  $k+1$ .

Proof:- The set  $C$  is known to both transmitter and the receiver, so if  $w \in W$  is the message and  $c = E(w)$  is transmitted, let  $c \notin T(c) = r$ . If the minimum distance between the code words is at least  $k+1$ , then the transmission of  $c$  ~~can~~ can result in as many as  $k$  errors and  $r$  will not be listed in  $C$ . Hence, we can detect all errors  $e$  where  $wt(e) \leq k$ .

(\*) If  $c_1$  and  $c_2$  are the codewords with  $d(c_1, c_2) < k+1$ . Then  $c_2 = c_1 + e$  where  $wt(e) \leq k$ . If we send  $c_1$  &  $T(c_1) = c_2$ , then we feel that if  $c_2$  is sent, then failing to detect an error of weight  $\leq k$ .

With  $E, W$ , and  $C$  are as above, and  $k \in \mathbb{Z}^+$ , we can construct a decoding function  $D: \mathbb{Z}_2^n \rightarrow W$  that corrects all transmission errors of weight  $\leq k$  iff minimum distance between code words is at least  $2k+1$ .

Proof: For  $c \in C$ , consider  $S(c, k) = \{x \in \mathbb{Z}_2^n \mid d(c, x) \leq k\}$ . Define  $D: \mathbb{Z}_2^n \rightarrow W$  as follows. If  $r \in \mathbb{Z}_2^n$  and  $r \in S(c, k)$  for some code word  $c$ , then  $D(r) = w$  where  $E(w) = c$ . If  $r \notin S(c, k)$  for any  $c \in C$ , then we define  $D(r) = w_0$ , where  $w_0$  is some arbitrary

The only problem we would face here is that  $D$  might not be a function. This will happen if there is an element  $r \in \mathbb{Z}_2^n$  with  $r$  in both  $S(g, k)$  and  $S(g_2, k)$  for distinct code words  $g, g_2$ . But  $r \in S(g, k) \Rightarrow d(r, g_1) \leq k$ , and  $r \in S(g_2, k) \Rightarrow d(r, g_2) \leq k$ , so  $d(g, g_2) \leq d(g, r) + d(r, g_2) \leq k+k = 2k$ . Consequently, if the minimum distance between code words is at least  $2k+1$ , then  $D$  is a function, and it will decode all possible received words, correcting any transmission error of weight  $\leq k$ . Conversely,  $g, g_2 \in C$  &  $d(g, g_2) \leq 2k$ , then  $g_2$  can be obtained from  $g$  by making at most  $2k$  changes. Starting at code word  $g$  we make approximately half of these changes. This brings us to  $r = g + e_1$  with  $\text{wt}(e_1) \leq k$ .

Continuing from  $r$ , we make the remaining changes to get to  $g_2$  and find  $r + e_2 = g_2$  with  $\text{wt}(e_2) \leq k$ . But then  $r = g_2 + e_2$ . Now with  $g + e_1 = r = g_2 + e_2$  and  $\text{wt}(e_1), \text{wt}(e_2) \leq k$ , how can one decide on the code word from which  $r$  arises? This ambiguity arises in a possible error of weight  $\leq k$  that cannot be corrected.

Ex: With  $W = \mathbb{Z}_2^2$  let  $E: W \rightarrow \mathbb{Z}_2^6$  given by (16)

$$E(00) = 000000, \quad E(10) = 101010, \quad E(01) = 010101, \quad E(11) = 111111.$$

$\Rightarrow$  The minimum distance between the code words is 3.

Hence, we can detect double errors & correct single ones.

$$S(000000, 1) = \{x \in \mathbb{Z}_2^6 \mid d(000000, x) \leq 1\}$$

$$= \{000000, 100000, 010000, 001000, 000100, 000010, 000001\}$$

The decoding function  $D: \mathbb{Z}_2^6 \rightarrow W$  gives

$$D(x) = \infty \text{ for all } x \in S(000000, 1)$$

lucy

$$\begin{aligned}
 S(010101, 1) &= \{x \in \mathbb{Z}_2^6 \mid d(010101, x) \leq 1\} \\
 &= \{010101, 110101, 011101, 010001, 010111, 010100\}
 \end{aligned}$$

here  $D(x) = 01$  for each  $x \in S(010101, 1)$

At this point our definition of D accounts for 14 of the elements in  $Z_2^6$ . Continuing to define D for the 14 elements in  $S(101010, 1) \cup S(111111, 1)$ , there remain 36 other elements to account for. We define  $D_W = \infty$  for these 36 elements and have a decoding rule that will correct single errors.

If  $c = 010101$  &  $T(c) = r = \boxed{0}1\boxed{0}101$ , we can detect this double error because  $r$  is not a code word. But if  $T(c) = r_1 = 111111$ , 3 errors, we think that  $c = 111111$ , & it decode to 11, (incorrectly). So it is count as 01.

Two bit errors can't be corrected.

## Parity check and Generator matrix :-

Here the encoding and decoding functions are given by matrices over  $\mathbb{Z}_2$ . One of these matrices will help us to locate the nearest code word for a received word. This is helpful when the set  $C$  grows larger.

Ex: Let  $G = \begin{bmatrix} 1 & 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 & 0 & 1 \end{bmatrix}$  be a  $3 \times 6$  matrix over  $\mathbb{Z}_2$ .

$$G = \left[ \begin{array}{c|c} I_3 & A \end{array} \right]$$

The first 3 columns of  $G$  constitute Identity matrix. The remaining 3 columns represent matrix  $A$ .

This partitioned matrix  $G$  is called generator matrix. Using this generator matrix  $G$ , we can define the encoding function  $E: \mathbb{Z}_2^3 \rightarrow \mathbb{Z}_2^6$  as follows.

$$\boxed{\text{For } w \in \mathbb{Z}_2^3, E(w) = wG}$$

**NOTE:** Assume that all of  $\mathbb{Z}_2^3$  is encoded and that the transmitter and the receiver will both know the real messages of importance and their corresponding code words.

for example,

$$E(110) = 110G = [1 \ 1 \ 0] \begin{bmatrix} 1 & 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 & 0 & 1 \end{bmatrix} = [110 \ 101]$$

$$\& E(010) = 010G = [0 \ 1 \ 0] \begin{bmatrix} 1 & 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 & 0 & 1 \end{bmatrix} = [010 \ 011]$$

Note that  $E(110)$  is obtained by adding first two rows of  $G$ .  $E(010)$  is simply the second row of  $G$ .

$W = \{000, 100, 010, 001, 110, 101, 011, 111\} \subseteq \mathbb{Z}_2^3$  and their code words are

$$C = \{000000, 100110, 010011, 001101, 110101, 101011, 011110, 111000\} \subseteq \mathbb{Z}_2^6$$

After receiving a codeword, one can recover the message by simply dropping the last 3 bits. The minimum

For all  $w = w_1 w_2 w_3 \in \mathbb{Z}_2^3$ ,

$$E(w) = w_1 w_2 w_3 w_4 w_5 w_6 \in \mathbb{Z}_2^6.$$

$$E(w) = [w_1 w_2 w_3] \begin{bmatrix} 1 & 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 & 0 & 1 \end{bmatrix}$$

$$= [w_1 w_2 w_3 (w_1 + w_3) (w_1 + w_2) (w_2 + w_3)]$$

so we have  $w_4 = w_1 + w_3$ ,  $w_5 = w_1 + w_2$ ,  $w_6 = w_2 + w_3$   
and these equations are called parity-check-equations.

Since  $w_i \in \mathbb{Z}_2$  for each  $1 \leq i \leq 6$ , it follows that  $w_i = -w_i$ ,  
and so these parity check equations are rewritten as

$$w_1 + w_3 + w_4 = 0$$

$$w_1 + w_2 + w_5 = 0$$

$$w_2 + w_3 + w_6 = 0$$

Thus we find that

$$\begin{bmatrix} 1 & 0 & 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \\ w_5 \\ w_6 \end{bmatrix} = H \cdot (E(w))^T = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

where  $(E(w))^T$  denotes the transpose of  $E(w)$ .

If the received word is  $r = r_1 r_2 \dots r_6 \in \mathbb{Z}_2^6$ , we can  
identify  $r$  as a code word iff

$$H \cdot r^T = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$H = [B \mid I_3] \quad \text{note that } B = A^T$$

In this example, since the minimum distance is 3, we are  
able to define a decoding function which corrects a  
single errors.

Suppose we receive  $r = 110110$ , we want to find the code word  $c$  which is nearest to  $r$ . Just compare  $r$  with each code allowable code word. If  $|C|$  is very high, just examine  $H \cdot r^T$ , this is called as syndrome of  $r$ . (18)

$$H \cdot r^T = \begin{bmatrix} 1 & 0 & 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}$$

so  $r$  is not a codeword, hence the error is detected.

Now compute the distance between  $r$  & each codeword  $c \in C$ , and consider the one which has less distance.

$$\therefore d(100110, r) = 1 \text{ & all other } c \in C, d(r, c) \geq 2.$$

Writing  $r = c + e = (100110 + 010000)$ , we find that the transmission error of (weight 1) occurs in the second component of  $r$ .

Hence, changing the second bit of  $r$ , we will get  $c$ , the message is comprised the first three components of  $c$ .

Let  $r = c + e$ , where  $c$  is a code word and  $e$  is an error pattern of weight 1. Suppose that 1 is the  $i$ th component of  $e$ , where  $1 \leq i \leq 6$ . Then

$$H \cdot r^T = H(c + e)^T = H \cdot (c^T + e^T) = H \cdot c^T + H \cdot e^T$$

With  $c$  the code word, it follows that  $H \cdot c^T = 0$ , so

$H \cdot r^T = H \cdot e^T = i$ th column of the matrix  $H$ . Hence  $c$  &  $r$  differ only in the  $i$ th component, and we can determine  $c$  by simply changing the  $i$ th component of  $r$ .

Suppose,  $r = 000111$ , its syndrome,

$$H \cdot r^T = \begin{bmatrix} 1 & 0 & 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

If  $c = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 1 \\ 1 \end{pmatrix}$  :  $\boxed{000111 = r}$