

Date / / DNA storage and security

CANVAS

→ IT495 - Course Code

Natural computing - bridge between computer science and natural sciences.

→ Mid term test - 30%
Scrib notes - 20%
Project - 50%

→ You should try to do I course in biology because computer science and biology combined can do wonders.

→ Eric Lander course - secrets of life (Beautiful course)

→ The next century is for biology because it can do wonders which we have never thought earlier.

→ "If you are doing some theory, you should know how to sell it."

→ Nature uses DNA for storage.

→ Genome: In molecular biology and genetics, genome is all genetic information of an organism. It consists of nucleotide sequences of DNA.

→ RNA computing, DNA computing, Bacterial computing
V. Imp.

Date / /

- In DNA,
 - A, T → complement
 - C, G → complement
(const form)
- ACCT is acronym for four types of bases found in DNA molecules:
 - A → adenine
 - C → cytosine
 - G → guanine
 - T → thymine
- cloud - bunch of computers distributed across networks
- storing the data is very much challenging task now, because data generated every day is very huge.
- for 1 Yottabyte storage, the cost is \$100 Trillion which is more than GDP of whole world.
- So, what is next generation storage device? DNA
- 1g of DNA can store 455 EB of data. (EB - Exabyte)
- Structure of DNA - search on YouTube
(Double helix structure)
- DNA synthesis - write process
- DNA sequencing - read process

Date ___ / ___ / ___

- Anything that is ever produced in this world can be stored in 1 KB of DNA.
- Access time for data stored on DNA is more than current storage technologies.
- Synthesizer machine - produces physical DNA
- Error correction is very very important in both reading from DNA and writing on DNA.
- Codec : A codec compresses or decompresses media files such as storage or videos.

⇒ Lecture 2 :-

→ Codec - Encoding & decoding method

→ DNA storage system steps :-

Encoding → synthesis → storage → retrieval → Search
encoding → decoding

(whole course is divided into above 6 parts)

→ Sequencing : reading of DNA

→ Error correcting codes are required for efficient storage and retrieval.

- Claude E. Shannon - father of Information Technology
(A mathematical Theory of Communication)
If you want to understand information theory, read this paper.
- Small molecules follow Clifford algebra.
- Regenerating codes technique is nowadays used in cloud computing
- Integer $x \cdot n$ is field if and only if n is prime number.
This is very important because the elements data packets are nothing but elements of finite fields.
- \mathbb{Z}_n is a field, if and only if n is a prime number.
 $n = \{2, 3, 4, 5, 6, 7, 8, 9, 10\}$
 $\times \quad \times \quad \times$

$$\text{GF}(4) = \{atb \mid a, b \in \mathbb{Z}_2, a^2 + at + t = 0\}$$

$$= \{0, 1, \alpha, \alpha^2\} \quad \{0, 1, \alpha, 1+\alpha\}$$

$$= \{0, 0, 1, 10, 11\}$$

(bit representation)

3 different representations

representations

We can construct field with any prime power.
In general we can construct,

$$\text{GF}(p^m) = \text{GF}(q)$$

where

q is prime power
 p is prime.

→ here we are talking about data which is discrete and that's why we require this discrete structure which is "finite field".

→ Search on youtube for finite fields.

$$\rightarrow \mathbb{Z}_2^n = \{(x_1, x_2, \dots, x_n) \mid x_i \in \mathbb{Z}_2\}$$

$$\mathbb{R}^n = \{(x_1, x_2, \dots, x_n) \mid x_i \in \mathbb{R}\}$$

$$n=3, \quad \mathbb{Z}_2^3 = \left\{ \begin{array}{l} 000 \\ 001 \\ 010 \\ 011 \\ 100 \\ 101 \\ 110 \\ 111 \end{array} \right\} \quad \begin{array}{l} \leftarrow \text{each vector in this} \\ \text{space is combination of} \\ \{000, 010, 100\} \end{array}$$

\uparrow
 \uparrow
 \uparrow

this is basis of
the vector space.

In \mathbb{R}^m , there are infinite vectors.

but \mathbb{Z}_2^n is finite.

hence basis for $\mathbb{Z}_2^3 = \{001, 010, 100\}$
dimension of $\mathbb{Z}_2^3 = 3 \leftarrow \text{size of basis vector}$

→ any nonzero vector forms a one dimensional space.

→ If basis $B = \{\bar{x}\}$, $\bar{x} \neq 0$ and $\bar{x} \in \mathbb{Z}_2^3$

$\bar{x} = \langle 110 \rangle$ space generated by 110

$$= \{\lambda(110)\} \quad \lambda \in \mathbb{Z}_2$$

$$= \{000, 110\}$$

Date _____

→ From out of q vectors we can form 7 one-dimensional service spaces, and 1 ~~the~~ 0 -dimensional service space.

→ How many 2 -dimensional service spaces are there?
you can choose any 2 which are linearly independent

→ 2 -dimensional service space - q vectors

$\begin{matrix} 1 & " & " & " \end{matrix}$ - 2 vectors

$\begin{matrix} 0 & " & " & " \end{matrix}$ - 1 vector

→ Search : what is dimension - $0, 1, 2$ service spaces?

→ $\mathbb{F}_q^n = GF(q)$ (field with q number of elements)

$$\mathbb{F}_q^n = \{ (x_1, x_2, \dots, x_n) \mid x_i \in \mathbb{F}_q \}$$

linear code is subspace of \mathbb{F}_q^n .

$$C \subseteq \mathbb{F}_q^n$$

✓

Linear code which is used for correcting errors is mathematically subspace of a vector space over finite field.

→ If vector space is of dimension $n \rightarrow$ the basis has n vectors

Date _____ / _____ / _____

⇒ lecture - 3→ code $C = \{000, 111\}$ is 1

here dimension of C because there is only 1 vector 111 which will generate whole space,

→ $\lambda(111)$ where $\lambda = 0$ or 1
so it generates whole space

→ $\mathbb{Z}_2^n \rightarrow$ vector space of dimension n
(linear code)

codespace $C \subseteq \mathbb{Z}_2^n$ k -dimensional

→ linear code C is simply a subspace of \mathbb{Z}_2^n .

In general ex: $C \subseteq \mathbb{Z}_2^3$ (3-dimensional)
 $\{000, 111\}$

→ Since it is k -dimensional service space it will have a basis
 $B = \{\bar{v}_1, \bar{v}_2, \dots, \bar{v}_k\}$

(1) all these k vectors are linearly independent

$$\sum_{i=1}^n \lambda_i \bar{v}_i = 0 \Rightarrow \lambda_i = 0 \text{ for } \forall i$$

(2) $\langle \bar{v}_1, \bar{v}_2, \dots, \bar{v}_k \rangle = C$

(Span of $\bar{v}_1, \bar{v}_2, \dots, \bar{v}_k$ is whole space C .)

$$\therefore \{ \lambda_1 \bar{v}_1 + \lambda_2 \bar{v}_2 + \dots + \lambda_k \bar{v}_k \mid \lambda_i \in \mathbb{Z}_2 \} = C$$

(λ_i are scalars in \mathbb{Z}_2)

Q) our code C which is used for correcting errors is nothing but a k -dimensional service space.

Date _____ / _____ / _____

so, total vectors will be 2^k , because A_2 has 2 choices : either 0 or 1.

→ for our example,

$$\{000, 111\} = \langle 111 \rangle = \{\lambda(11) \mid \lambda \in \mathbb{Z}_2^*\}$$

We can write it in the form of matrix:

$$Q_L = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_k \end{bmatrix}_{k \times n}$$

Hence in above example Q_1 is 111
 $\therefore Q_1 = [111]$

→ how encoding is done?

$$(x_1, x_2, \dots, x_k) G_{k,m} = (y_1, y_2, \dots, y_m)$$

if bit is 0 \Rightarrow 0 [111] = [000]

If bit is 1 \Rightarrow 1[11] = [11]

90, 0 → 000

1 - 111

$$C = \begin{Bmatrix} 0 & 0 & 0 \\ 1 & 1 & 1 \end{Bmatrix}$$

So, error correcting codes are subspace of vector space.

Date _____ / _____ / _____

if the Hamming distance is d , you can correct $\frac{(d-1)}{2}$ errors

→ encoding is easy, decoding is difficult.

→ $C \subseteq \mathbb{Z}_2^n$

This linear code is described as,

$O : [n, k, d]_2$ ($_2$ means binary)

$n \rightarrow$ length of \in (subspace)

$k \rightarrow$ dimension

$d \rightarrow$ Hamming distance

minimum

If d is large we can correct many errors.

$$d = \min \{ d_K(\bar{x}, \bar{y}) \mid \bar{x}, \bar{y} \in C \}$$

example

$$A = \begin{bmatrix} 110 \\ 101 \\ 011 \end{bmatrix}_{2 \times 3} \rightarrow n=3$$

$\downarrow k=2$

$$\Rightarrow C = \{000, 110, 101, 011\}$$

$$\text{so, } C : [3, 2, 2]_2$$

so, we can say it is binary code of length 3, dimension 2 and distance 2.

this is not useful because $(d-1) = 0$ so, you can correct 0 errors. d should be at least $\frac{3}{2}$ so that you can correct errors

Date / /

$$\rightarrow (01) \begin{bmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} = [101]$$

so, (01) is encoded as 101

so, n is important which is dimension of service space and n information bits are converted to n bits.

\rightarrow Dual space : All those vectors in a vector space such that their inner product is 0.

C^\perp
↑
symbol

$$C^\perp = \{ \bar{y} \in \mathbb{Z}_2^n \mid \langle \bar{x}, \bar{y} \rangle = 0, \forall \bar{x} \in C \}$$

$$\text{where } \langle \bar{x}, \bar{y} \rangle = x_1 y_1 + x_2 y_2 + \dots + x_n y_n = 0$$

if $n=3$ and $C=\{000, 111\}$, then,

$$\mathbb{Z}_2^3 = \{ \begin{array}{l} 000 \\ 001 \\ 100 \\ 010 \\ 110 \\ 101 \\ 011 \\ 111 \end{array} \}$$

so, C^\perp is set of all vectors in \mathbb{Z}_2^3 such that its inner product with all vectors in C is 0.

$$C^\perp = \{ 000, 110, 101, 011 \}$$

$$\text{because } \begin{aligned} \langle 110, 000 \rangle &= (0+0+0) \cdot 2^0 = 0 \\ \langle 110, 111 \rangle &= (1+1+0) \cdot 2^1 = 0 \end{aligned}$$

Date / /

What is dimension of C^\perp ?

ans: 2

Remember that,

$$\dim C + \dim C^\perp = \dim \mathbb{Z}_2^3$$

$$\therefore 1 + 2 = 3$$

(Note: Here we are ignoring 000 in calculating dimension)

- (1) C is self orthogonal if $C \subseteq C^\perp$
- (2) C is self dual if $C = C^\perp$

$$\rightarrow G = \begin{bmatrix} I_k & | & A \\ \hline & I_{n-k} \end{bmatrix}_{K \times n} \quad (A \text{ is generator matrix})$$

$H_{n \times n}$ is parity check matrix ($P(C)$) and its size is $(n-k) \times n$

$$H = [-A^T \mid I_{n-k}]$$

Note: H is matrix such that $GHT^T = 0$

→ example

$$G = \begin{bmatrix} 1 & 0 & 1 & 1 \\ 0 & 1 & 0 & 2 \end{bmatrix}_{K \times n} = \begin{bmatrix} I_2 & | & A \end{bmatrix}_{2 \times 4} \quad K=2 \quad n=4$$

$$\Rightarrow A = \begin{bmatrix} 1 & 1 \\ 0 & 2 \end{bmatrix} \Rightarrow H = [-A^T \mid I_{n-k}]$$

$$H = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 \end{bmatrix} \rightarrow \text{Parity check matrix}$$

Date / /

So, any subspace can be described in two ways. (R matrices)

(1) Generatormatrix (G)

(2) Parity check Matrix (H)

$$C = \langle G \rangle$$

$$= \langle \bar{v}_1, \bar{v}_2, \bar{v}_m \rangle$$

$$C = \{ \bar{x} \in \mathbb{Z}_2^n \mid H \bar{x}^t = 0 \}$$

so, many times you describe subspace by describing H.
C is null space of H. i.e. all those vectors \bar{x} such
that $H \bar{x}^t = 0$

→ lecture - 4

Date: 25/06/2022

→ code is collection of binary strings of some length. code is basically subset of \mathbb{Z}_2^n .

→ code is called linear if it is subspace. for example:
 $C = \{000, 111\}$ is subspace of \mathbb{Z}_2^3 .

→ $\mathbb{Z}_2^n = \{ (x_1, x_2, \dots, x_n) \mid x_i \in \mathbb{Z}_2 \}$

(Code)

$C \subseteq \mathbb{Z}_2^n$

subset

$C : (n, m \geq d)$

length

minimum Hamming distance

$|C| \leftarrow \text{Size of } C$

= # of codewords

= dimension

so, code is simply subset of whole space. You can do encoding, decoding for these kind of codes and these are called nonlinear codes. Nonlinear means if you take two vectors

out of the set C, their sum may not belong to same set C.

and linear code is simply k -dimensional generic space.

see the left-in $\hookrightarrow G \leq Z_2^n$

Symbol k-dim Subspace

$C: [n, k, d]$

length ↑ dimension ↑ ↑

Minimum Hamming Distance

(H17P) For ex: $C = \{000, 111\}$ is $\begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$ code

length \rightarrow \uparrow \uparrow \leftarrow min Hamming distance

dim (because there is only 1 vector in basis)

c can be described by 2 kind of matrices,

(1) Generator matrix (G_0)

(Q) Parity check matrix (H)

Method -1

generator matrix : $G_{km} = [I_k | A]$ (standard form of generator matrix)

When you write down all the basis vectors of C in matrix form you will get generator matrix $G_{n \times m}$.

→ as there are k -linearly independent vectors in basis of C , there will be in total 2^k vectors in C . The k -linearly independent vectors span whitespace.

$\rightarrow c$ is linear combination of n -linearly independent vectors.

\rightarrow cardinality of E is 2^K .

Method-2 : Parity check Matrix

$H_{(n-k) \times n}$ P.C.M. binary matrix

$$C = \{ \bar{x} \in \mathbb{Z}_2^n \mid H\bar{x}^T = 0 \}$$

→ null space of matrix H

If, $\alpha = [I_k \mid A]$ then, $(I - \text{box})$
 $H = [-A^T \mid I_{n-k}]$ $n-k \times n$

→ Rank of α ? (How many L.I. rows or columns?)

ans : k

→ Rank of H? ans : $n-k$

→ Relationship between α and H is,

$$\alpha H^T = 0$$

example :

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

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

\checkmark
 $C = \{0000, 1011, 0101, 1110\}$ (by taking linear combination from α)

also, $C = \{ \bar{x} \in \mathbb{Z}_2^n \mid H\bar{x}^T = 0 \}$ (we want to verify this).
 Here n is 4

$$H\bar{x}^T = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = 0 \Rightarrow \begin{array}{l} x_1 + x_3 = 0 \\ x_1 + x_2 + x_4 = 0 \end{array}$$

remember we want to encode k information bits to n bits

So, matrix H works like,

we have k info. bits, we extend it to n info. bits by adding $(n-k)$ bits more such that $H\vec{x}^t = 0$ for some matrix H .

So, we are adding parity check.

which is $(n \times k \times n)$

here $k=2$ $n=4$

$$\rightarrow x_1 + x_3 = 0$$

$$\rightarrow x_3 = -x_1 = x_1$$

$$\rightarrow x_1 + x_2 + x_4 = 0$$

$$\Rightarrow x_4 = -x_1 - x_2 = x_1 + x_2$$

x_1, x_2 we can control and from that we can get x_3 and x_4 .

| x_1 | x_2 | x_3 | x_4 |
|-------|-------|-------|-------|
| 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 1 |
| 1 | 0 | 1 | 1 |
| 1 | 1 | 1 | 0 |

you can see we are getting same vectors as in C.

vectors of length k

If we are using generator matrix (G), multiply with G which is $k \times n$ and we get n bits.

vector of length

If we have parity check matrix, we add extra $n-k$ bits such that $H\vec{x}^t = 0$.

Date / /

- if minimum Hamming distance is d , you can correct $(\frac{d-1}{2})$ errors.
- for $d=2$, you may be able to correct some errors but not all errors.

IMP

$$\rightarrow \text{syn}(y) = Hy^t \quad (\text{syndrome decoding})$$

syndrome

This is used to simplify decoding procedure
you can compute syndrome for each vector.

→ (See in slides the value of $\text{syn}(y)$ for vectors)

→ y is received bits

→ for decoding using syndrome decoding,
if x is sent and y is received, just compute the
syndrome,

$$\text{syn}(y) = Hy^t$$

and whichever bit is 1 in $\text{syn}(y)$, flip that bit
in received sequence.

Search : What is syndrome decoding?

$$C = \{000, 111\} \subset \mathbb{Z}_2^3$$

find the sphere of radius t along vector y

$$S_t(\bar{x}) = \{ \bar{y} \in \mathbb{Z}_2^3 \mid d(\bar{y}, \bar{x}) \leq t \}$$

$$t = \left[\frac{d-1}{2} \right]$$

Date / /

$$S_1(000) = \{ \quad \}$$

$$S_2(111) = \{ \quad \}$$

decoding is NP-complete problem.

Lecture 5 :-

Date: 25/08/2022

Z_2^m consists 2^m vectors.

for $m=2$, $Z_2^2 = \{ 00, 01, 10, 11 \}$

$$H = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \end{bmatrix}_{2 \times 3} \quad (\text{all vectors other than } 00)$$

for $m=3$, $Z_2^3 = \{ 000, 011, 101, 110, 111, 100, 010, 001 \}$

$$H = \begin{bmatrix} 0 & 1 & 1 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 & 1 \end{bmatrix}_{3 \times 7} \quad (\text{just write binary expression of numbers other than 0 as columns})$$

(all vectors other than 000 are used as columns)

→ make sure that no 2 columns are multiple of each other.

→ this parity check matrix generates the code known as 'Hamming code'.

Parameters for this code is $[7, 4, 3]$

$\begin{matrix} n \\ k \\ d \end{matrix}$

\bar{x} , here - indicates that it is vector.

Date / /

$$G = \{ \bar{x} \in \mathbb{Z}_2^3 \mid H\bar{x}^t = 0 \}$$

(HHP) → How many vectors will be there in C ?
 ans $= 2^4$ because dimension is $k=4$. so total possible
 vectors are $2^4 = 16$.

$k=4$ indicates there are 4 vectors in generator matrix.

In general we can construct code of length $= 2^{m-1}$, dimension $= 2^{m-1-m}$ and distance = 3.

$$\therefore H_m : [2^{m-1}, 2^{m-1-m}, \beta],$$

$$H_m : [2^{m-1}, 2^{m-1-m}, \beta]_2$$

$$\downarrow \quad \downarrow \quad \downarrow$$

$$n \quad k \quad d$$

since distance is 3, it can correct $(d-1)/2 = 3/2 = 1$ error.

→ Decoding

If $\bar{x} \in G$ is sent and \bar{y} is received, then,

$$\bar{y} = \bar{x} + \bar{e}$$

$$\bar{e} = (\epsilon_1 \epsilon_2 \dots \epsilon_n)$$

(error vector)

$$\rightarrow \text{Syn}(\bar{y}) = \text{Syn}(\bar{x} + \bar{e})$$

$$= H\bar{x}^t + H\bar{e}^t$$

but we know that if $\bar{x} \notin G$ then $H\bar{x}^t = 0$

Date _____ / _____ / _____

$$\therefore \text{syn}(\bar{y}) = H\bar{e}^t$$

let $H = (H_1 \ H_2 \ \dots \ H_n)$ and $e = \begin{pmatrix} 0 & 0 & \dots & 0 & 1 & 0 & 0 & 1 & \dots & 0 & 0 & 0 \end{pmatrix}$
 $a \ b \ c$

a, b, c are positions on which 1 is occurring in e vector

$$\text{syn}(\bar{y}) = \sum_{i=1}^n e_i H_i = H_a + H_b + H_c$$

so, syndrome is basically sum of columns in ~~with~~ parity check matrix where error occurred.

$$H = \begin{bmatrix} 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ 0 & 1 & 1 & 0 & 0 & 1 & 1 \\ 1 & 0 & 2 & 0 & 1 & 0 & 1 \end{bmatrix}_{3 \times 7}$$

$$H: [7 \ 4 \ 3]$$

\uparrow \uparrow \uparrow
 n k d
 (length) (dimension)

suppose received vector is $\bar{y} = 1 \ 1 \ 0 \ 1 \ 0 \ 1 \ 1$

$$\begin{aligned} \text{syn}(\bar{y}) &= H\bar{y}^t \\ &= \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} \quad (6\text{th column}) \end{aligned}$$

this means that 6th bit is corrupted.

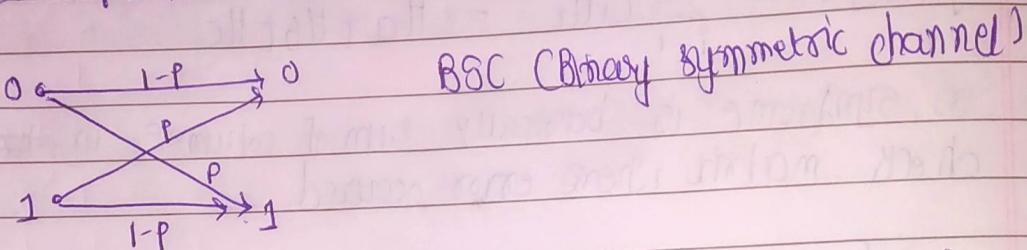
so, $\bar{x} = 1 \ 1 \ 0 \ 1 \ 0 \ 0 \ 1$ (we reversed 6th bit)

Date _____

→ weight of the code is $\frac{k}{n} \rightarrow \frac{\text{dimension}}{\text{length}}$

→ see [codetables.de](http://www.codetables.de) website

→ There are only 3 examples of code which are posted
 (1) repetition code (2) Hamming code (3) Golay code



see the example of probabilistic decoding in ppt.

→ coding theory and cryptography are connected, we can say coding theory is sister of cryptography

⇒ cyclic code

Properties :-

(1) Linear code (there should be basis vector)

(2) any cyclic shift of codeword is again a codeword
 i.e. $c_0 c_1 \dots c_{n-1} \in C \Rightarrow c_{n-1} c_0 \dots c_{n-2} \in C$

ex: 1 {000, 101, 011, 110}

this is cyclic code because it is linear code and also it follows 2nd property,

⇒ {0000, 1001, 0110, 1111} ← ex: 2

→ is not cyclic because 2nd property is not followed

Date / /

$$\Rightarrow \{0000, 1010, 0101, 1111\} \xleftarrow{\text{ex:3}}$$

this is cyclic code because both the properties of cyclic codes are followed.

$$\Rightarrow \text{ex:4 } \{11111\}$$

not cyclic code because it is not linear because
00000 is not there

$$\Rightarrow \{\bar{0}\} = \{000\dots 0\}$$

this is not cyclic code.

$$\Rightarrow \{\lambda I, \lambda \in \mathbb{F}_q\}$$

It is cyclic code.

$$\Rightarrow \mathbb{F}_q^n \quad (n\text{-dimensional vector space over } q \text{ elements})$$

It is also cyclic code.

Mapping to polynomial for representing cyclic codes

$$c_0 c_1 \dots c_{n-1} \longleftrightarrow c_0 + c_1 x + c_2 x^2 + \dots + c_{n-1} x^{n-1}$$

$c(x)$

$$c(x) = c_0 + c_1 x + c_2 x^2 + c_3 x^3 + \dots + c_{n-1} x^{n-1}$$

multiplied by x^n
means 1 cyclic shift

now x^n is going outside of space because

degree of polynomial is going outside

$$\Rightarrow x^n = 1$$

So, remember that whenever in polynomial you encounter x^n , it is 1.

Lecture 6

Z_2^3 $Z_2[\alpha]$ $\alpha + \alpha x + \alpha x^2 + \dots + \alpha x^{n-1}$

$$000 = 0$$

$$100 = 1$$

$$010 = \alpha$$

$$001 = \alpha^2$$

$$110 = 1+\alpha$$

$$011 = \alpha+\alpha^2$$

$$101 = 1+\alpha^2$$

$$111 = 1+\alpha+\alpha^2$$

$$\rightarrow f(x) = x^2 + x + 1$$

$$\rightarrow Z_2[\alpha] \xrightarrow{\text{Modulo}} (x^2 + x + 1) \rightarrow GF(4)$$

$$\alpha^2 = \alpha + 1 \quad (\alpha^2 + \alpha + 1 = 0)$$

| | | | | |
|--------------|--------------|--------------|--------------|--------------|
| $+$ | 1 | α | $\alpha + 1$ | 0 |
| 1 | 0 | $\alpha + 1$ | α | 1 |
| α | $\alpha + 1$ | 0 | 1 | α |
| $\alpha + 1$ | α | 1 | 0 | $\alpha + 1$ |
| 0 | 1 | α | $\alpha + 1$ | 0 |

$$\rightarrow Z_2[\alpha] / (x^2 + 1)$$

Modulo \nearrow

| | | | |
|------------|------------|------------|------------|
| * | 1 | α | $1+\alpha$ |
| 1 | 1 | α | $1+\alpha$ |
| α | α | 1 | $1+\alpha$ |
| $1+\alpha$ | $1+\alpha$ | $1+\alpha$ | 0 |

Hence $(1+\alpha)$ does not have inverse. So it is not a field.
(there is no column as 1 for $(1+\alpha)$.)

$$\rightarrow \mathbb{Z}_2[\alpha] / \langle \alpha^3 - 1 \rangle = R_3(\text{Ring})$$

here $n=3$ so, $\alpha^3 = 1$ (should be 1)

$$\boxed{\mathbb{Z}_2[\alpha] / \langle \alpha^3 - 1 \rangle = R_3(\text{Ring})}$$

in general,

$$\mathbb{Z}_2[\alpha] / \langle \alpha^n - 1 \rangle = R_n(\text{polynomial ring})$$

→ let's see what happens for subspace.

$$C = \left\{ \begin{array}{l} 000 \\ 110 \\ 101 \\ 011 \end{array} \right\} \subset \mathbb{Z}_2^3$$

$$\text{Im } C = \{0, 1+\alpha, 1+\alpha^2, \alpha+\alpha^2\} \subset R_3 \text{ ideal}$$

(it is called ideal code for R_3)

Ring $I \subseteq R$ is ideal if
 $(\neq \emptyset)$

- (1) $a+b, a-b \in I$, $\forall a, b \in R$
 (2) $r.a \in I$, $\forall r \in R$
 $\forall a \in I$

→ For cyclic code just find Ideal subset.

$$R_3 = \mathbb{Z}_2[\alpha] (\alpha^3 - 1)$$

$$I = \{1, 1+\alpha, 1+\alpha^2, \alpha+\alpha^2\}$$

$$I = \langle 1+\alpha \rangle$$

$$\forall (1+\alpha)^k, \forall k \in \mathbb{Z}$$

→ I is prime ideal if $\exists a \in I$

$$\text{s.t. } I = \langle a \rangle$$

$$\{rg \mid r \in R\}$$

→ $G (\neq 0)$ cyclic code in $R_m = \mathbb{Z}_2[\alpha]/\langle \alpha^n - 1 \rangle$

(2) $\exists a$! monic polynomial $g(x)$ of smallest degree in $\mathbb{Z}_2[x]$

$$(x) - C = \langle g(x) \rangle$$

(3) $g(x)$ is factor of $(\alpha^n - 1)$

Date _____

→ If $g(x) = g_0 + g_1x + \dots + g_r x^r$ then,
 C is generated by

$$G = \begin{bmatrix} g_0 & g_1 & g_2 & g_3 & \dots & 0 & \dots & 0 \\ 0 & g_0 & g_1 & g_2 & \dots & g_3 & \dots & 0 \\ 0 & 0 & \dots & 0 & g_0 & g_1 & \dots & g_3 \end{bmatrix} \leftarrow \underline{\text{IMP}}$$

How to construct
 G from $g(x)$

→ How to find all binary cyclic codes of length 3?

$$x^3 - 1 = (x+1)(x^2 + x + 1)$$

Generator poly

code in \mathbb{R}_3 code in \mathbb{Z}_2^3

1

All of \mathbb{R}_3 all of \mathbb{Z}_2^3 $x+1$ {0, 1+x, 1+x^2, $x+x^2$ }

{000, 110, 011, 101}

 x^2+x+1

{0, 1+x+x^2}

{000, 111}

 x^3-1

{0}

{000}

→ n factorization of $x^n - 1$, no of cyclic codes.

CANVAS

Date _____ / _____ / _____

$(H^w) \rightarrow$

for $n=7$ which polynomial should we use for cyclic codes

Some interesting cyclic codes

Read solomon code

Read muller code

Hamming code

BCH code

mother of all codes

$$g(x) =$$

$$g(x) =$$

$$g(x) =$$

$$g(x) =$$

⇒ Lecture 7

Binary string - 0, 1

ternary || - 0, 1, 2

quaternary || - 0, 1, 2, 3

Most common error is deletion. chunk of letters are deleted.

→ If the encoded string has high minimum distance, we can correct more and more errors.

→ Edit distance: minimum number of insertion, deletion or substitution required to convert one string to another.

Goldman, church papers.

CD, DVD, QR code, High definition TV, deep space communication - uses Reed Solomon code

$m_0 + m_1x + m_2x^2 + \dots + m_{k-1}x^{k-1}$ - degree k-1 polynomial from message bits.

$\Sigma_{DNA}(n, n, d) \subseteq \Sigma_{DNA}^n$

DNA code is subset of Σ_{DNA}^n with length n and has

Date ___ / ___ / ___

m code words with minimum hamming distance d.

→ Main objective :- optimal DNA code with maximum minimum distance d (to correct more errors)

- Maximum storage capacity achieving codes.
- Designing the chunk architecture of DNA blocks.

→ Run length Encoding :

- lossless method
- text or programs
- Simplest method of compression

Ex: I/P : AAA BBB CCC DDD

Encoded : 3A 2B 1C 1D

→ Non Homopolymorphes mapping

| previous not written | next trait to encode | | | → using this map, the corresponding DNA string has |
|-------------------------|----------------------|---|---|--|
| | 0 | 1 | 2 | |
| A | C | G | T | |
| G | G | T | A | |
| C | T | A | C | no repetition of A, C, G, T. |
| T | A | C | G | |

AA, CC, GG, TT never occur. (non Homopolymorphes)

Date ___ / ___ / ___

→ 3 perfect codes,

(1) repetition code

(2) Hamming code

(3) Golay code

→ Generator matrix for extended Golay code

$$[84, 22, 8]_2$$

total 2^{12} codewords.

→ How ternary Golay code can be constructed through quadratic residue?

→ $[11, 6, 5]_3$ - Ternary Golay code

$$\# \text{ of codewords} = 3^6 = 729$$

⇒ lecture - 6

code rate = $\frac{k}{n}$, k bits are encoded to n bits

challenge is to construct a code which has code rate close to 1.

flash memory

→ CD, DVD, optical storage devices, Blue ray disks, use constrained coding

→ constrained coding: input is restricted.

Date / /

- The # of consecutive 0's and 1's is called as runlength. for certain applications, the runlength value lies between some min value and maximum value.
- suppose $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_n) \in \{0,1\}^n$

A dk-sequence has 2 constraints.

d constraint : consecutive 1's are separated by at least d zeroes.

1 1 1 ... 1 0 0 ... 0 1 1 1 ... 1

$\downarrow d$ zeroes (d zeroes)

$\downarrow k$ zeroes (at most k zeroes)

K constraint : consecutive 1's are separated by at most k zeroes.

→ RLL (run length limited sequence)

$$\bar{\alpha} = (\bar{\alpha}_1, \bar{\alpha}_2, \dots, \bar{\alpha}_n)$$

$$\bar{\alpha}_i = \bar{\alpha}_{i-1} \oplus \alpha_i, \bar{\alpha}_0 = 1$$

here α is dk-sequence

→ so from dk-sequence we are constructing RLL sequence.

$$\bar{\alpha} = 010010001$$

$$\bar{\alpha} = 100011110$$

Run length of all sequence lies between d_{t2} and k_{t1}

Date / /

→ capacity of constrained channel is,

$$C = \lim_{T \rightarrow \infty} \frac{1}{T} \log_2 N(T)$$

$N(T) = \# \text{ of allowed signal sequences in the time interval } T.$

→ $N_d(n) = \# \text{ of } d \text{ sequence of length } n$
 (K constraint is absent)

$$N_d(n) \approx C \lambda^n \quad n \geq 1$$

C is constant and λ is growth factor
 $1 < \lambda < e$, is the largest real root of characteristic equation.

$$z^{d+1} - z^d - 1 = 0$$

for such d constrained channel,
 channel capacity $[C_d = \log_2 \lambda]$

→ code efficiency = code rate $\frac{d}{\text{capacity of the constrained channel having same runlength constraint}}$

⇒ lecture 9

net information density = $\frac{\text{number of input information bits}}{\text{number of bases in resulting DNA storing}}$

Date / /

$$\rightarrow \text{GF}(q) = \{0, 1, \alpha_1, \dots, \alpha_{q-2}\} \quad q = p^m$$

T

 \mathbb{Z}_p for $\beta \neq 0 \in \text{GF}(q)$,order of β : smallest positive integer m such that $\beta^m = 1$

$$\text{ex: } \mathbb{Z}_5 = \{0, 1, 2, 3, 4\}$$

$$1^1 = 1 \Rightarrow \text{order}(1) = 1$$

~~$2^3 = 3$~~

$$2^4 = 1 \Rightarrow \text{order}(2) = 4$$

Similarly, $\text{order}(3) = 4$ and $\text{order}(4) = 2$

\rightarrow hence 3 and 2 are special because they are having
~~order~~ $= 5-1=4$. Basically they are generating
 'whole \mathbb{Z}_5 except 0'.

So, 3 & 2 are known as "primitive elements".

\rightarrow If you take any element in the field, its order
 will divide the order of multiplicative group (i.e. $(q-1)$)

\rightarrow each finite field has one or more elements
 which are PRIMITIVE ELEMENTS.

\rightarrow The # of elements in $\text{GF}(q)$ of order t is
 $\phi(t)$

Date / /

$\phi(t)$ is Euler totient function.

$$\phi(10) = 4$$

$$\phi(6) = \phi(2 * 3) = \phi(2) * \phi(3) = 2$$

$$\phi(15) = 8$$

→ GF(q): # of primitive elements is exactly $\phi(q-1)$

→ GF(2^4) = {0, 1, α , α^2 , ..., α^{14} }
 where α is root of polynomial,
 $\pi(\alpha) = 1 + \alpha + \alpha^4$.
 $\Rightarrow \alpha^4 + \alpha + 1 = 0$

α , α^2 , α^4 are primitive element whereas α^3, α^5 are not.

→ A primitive polynomial (P.P) is a polynomial having primitive element as its root.

→ finite fields - 2009 - 42 pg

→ If α is primitive element then α^t is also primitive

$$\alpha^{-1} = \alpha^{7-1} = \alpha^6$$

$$\alpha^{-2} = \alpha^{7-2} = \alpha^5$$

$$\alpha^{-4} = \alpha^{7-4} = \alpha^3$$

Date / /

If $p(\alpha)$ is primitive, then $p^*(\alpha) = \alpha^m p(\alpha^{-1})$

→ Number of primitive polynomials of degree m is $\frac{\phi(q-1)}{m}$.

→ bch codes - 1009, 4 pg

→ Set $\mathbb{Z}_{p^m-1} \times P$
this divides \mathbb{Z}_{p^m-1} into cosets.

cyclotomic coset containing g is,

$$C_0 = \{g, pg, p^2g, p^3g, \dots, p^{m_g-1}g\}$$

m_g is the smallest +ve integer such that

$$p^{m_g} \cdot g \equiv g \pmod{p^m-1}$$

Example

$$\text{mod } 15, \quad p=3 \quad (\mathbb{Z}_{15})$$

$$C_0 = \{0\}$$

$$C_1 = \{1, 2, 4, 8\}$$

$$C_3 = \{3, 6, 12, 9\}$$

$$C_5 = \{5, 10\}$$

$$C_7 = \{7, 14, 13, 11\}$$

These are all cyclotomic cosets.

→ lecture 10

Reed-Solomon codes

Applications: Flash drives, CDs, DVDs, QR code etc

Date _____

→ 8 bits correspond to element of a field with 256 elements.

→ The integer representation is divided into blocks of size K .

→ $\text{GF}(q)$

$$\text{message} : (m_0 m_1 m_2 \dots m_{K-1}) \in (\text{GF}(q))^K$$

$m_i \in \mathbb{Z}_2$

$$\text{polynomial } p(x) = m_0 + m_1 x + m_2 x^2 + \dots + m_{K-1} x^{K-1}$$

A Reed Solomon codeword $\bar{c} \in \text{RS}$

$$\left\{ \bar{c} = c_0 c_1 \dots c_{q-1} = p(0) p(\alpha) p(\alpha^2) \dots p(\alpha^{q-1}) \right\}$$

$c_i = p(\alpha^i)$ and

$c_0 = p(0)$ ↙

$$\rightarrow \text{GF}(q) = \{0, 1, \alpha, \alpha^2, \dots, \alpha^{q-2}\}$$

$$|\text{RS}| = q^K \quad (\text{size of Reed Solomon codes})$$

and RS codes are "linear codes"

$$\dim(\text{RS}) = K$$

$$n = q \quad (\text{length of each codeword}), \quad [n \xrightarrow{q} K \& d]$$

$$p(0) = m_0$$

$$p(\alpha) = m_0 + m_1 \alpha + \dots + m_{K-1} \alpha^{K-1}$$

$$p(\alpha^2) = m_0 + m_1 \alpha^2 + \dots + m_{K-1} \alpha^{2(K-1)}$$

Date / /

$$p(\alpha^{q-1}) = m_0 + m_1 \alpha^{q-1} + \dots + m_{K-1} \alpha^{(q-1)(K-1)}$$

Any K of these equations,

$$(X) \quad \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 1 & \alpha & \alpha^2 & \dots & \alpha^{K-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & \alpha^{K-1} & \alpha^{2(K-1)} & \dots & \alpha^{(K+1)(K-1)} \end{bmatrix} \begin{bmatrix} m_0 \\ m_1 \\ \vdots \\ m_{K-1} \end{bmatrix} = \begin{bmatrix} p(0) \\ p(\alpha) \\ \vdots \\ p(\alpha^{K-1}) \end{bmatrix}$$

→ vandermonde matrix is non-singular matrix as long as α_i are different.

→ here (X) is a nonsingular matrix. (out of q , t bits are corrupted)

→ suppose that t of the codeword's coordinates are corrupted.

we might construct all possible distinct system of K equations. So there are total q^K such systems. and $\binom{t+K-1}{K}$ of this will give incorrect solution.

→ Majority among all the solutions will give correct solution as long as,

$$\binom{t+K-1}{K} < \binom{q-1}{K} \quad \text{Why?}$$

$$\Rightarrow t+K-1 < q-1$$

$$\Rightarrow 2t < q-K+1 \quad \Rightarrow t = \left\lceil \frac{q-K+1}{2} \right\rceil$$

→ Reed-Solomon code has property that distance d is

$$d = n - k + t$$

Condition for
whether it is RS code
is $n \geq k+t$

$$\Rightarrow [n, k, n-k+t] \xrightarrow{d} \text{MDS}$$

(Maximum Distance Separable Codes)

You can find d if t is known.

→ If $m(x) = m_0 + m_1(x) + \dots + m_{k-1}x^{k-1}$

$$(Ex) c(x) = m(x)g(x) = \langle g(x) \rangle$$

$g(x)$ is generator polynomial.

$$g(x) = \prod_{j=1}^{q-1} (x - \alpha^j) \quad n = q-1$$

→ A polynomial having primitive elements as root is "primitive polynomial". primitive elements has order $(q-1)$. (q is order of field)

→ Reed-Solomon codes → special case of BCH codes.

⇒ Lecture-21

PCR : polymerase chain reaction (Makes copy of DNA)

DNA sequencer : machine for reading DNA

Transcription

Date / /

→ syndrome vector: $\text{Sym}(\bar{y}) = H\bar{y}^t = \begin{pmatrix} s_1 \\ s_2 \end{pmatrix}$

H is parity check matrix and \bar{y} is received vector.

- Hamming code is cyclic code and we can extend it to BCH code which can correct 2 errors.
- If H is parity check matrix and d is minimum distance of the code then any $(d-1)$ columns of H are linearly independent.
- $m_{i, \alpha}$ is minimal polynomial of α^i .
- Minimal polynomial is a polynomial of lowest degree having α as its root.

primitive polynomial is a polynomial which is having its root as primitive elements.

primitive polynomial can be minimal polynomial also. but every minimal polynomial may not be primitive polynomial.