

CS250/EE387 - LECTURE 6 - MAKING RS CODES BINARY

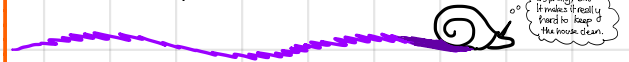
AGENDA

- ① BCH Codes
- ② Reed-Muller Codes [part I]
- ③ Concatenated Codes [part I]

The story so far is:

GASTROPOD FACT.

Some snails produce a secretion that is a color-fast natural dye. For example, "Tyrian purple" as well as other purple and blue dyes, were made from snail secretions.



Reed-Solomon Codes have the optimal trade-off between rate and distance, and they have efficient decoding algorithms!

That's a pretty good story. (Maybe we should just stop here?)

HOWEVER, there are some downsides to RS codes. A big one is the alphabet size.

GOAL. Obtain EXPLICIT (aka, efficiently constructible), ASYMPTOTICALLY GOOD families of BINARY CODES, ideally with fast algorithms.

Today we will see a few ways to approach this problem.

The first several won't work, but they are independently interesting and will come back later.

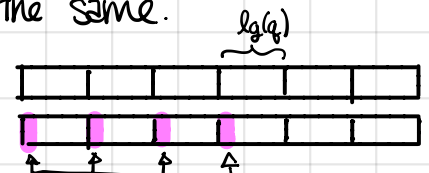
STRAWMAN. Replace every symbol of \mathbb{F}_q with $\log_2(q)$ bits.

If we start w/ n', k' RS code:

NOTE: This is actually done in practice!

$$\text{Rate} = \frac{k' \cdot \log_2(q)}{n' \cdot \log_2(q)} = \frac{k'}{n'} \text{ stays the same.}$$

$$\text{Relative Distance} = \frac{n' - k' + 1}{n' \cdot \log_2(q)} \sim \frac{1 - R}{\log(n')} \text{ if } n' = q - 1$$



corrupting d bits also can corrupt d symbols. So the distance is $d \geq n' - k' + 1$

$$\sim \frac{1 - R}{\log\left(\frac{n}{\log(n)}\right)} \sim \frac{1 - R}{\log(n)}$$

where $n = n' \log(n')$ is the block length of the binary code.

So the STRAWMAN is NOT asymptotically good.
 If R is constant, then $\delta \rightarrow 0$.

① BCH Codes. BCH = Bose and Ray-Chaudhuri, Hocquenghem

What if we just take $RS(n, k) \cap \{0, 1\}^n$?

DEF. Let $n = 2^m - 1$, let γ be a primitive element of \mathbb{F}_2^m .
 Then for $d \leq n$, define

$$BCH(n, d) = \left\{ (c_0, \dots, c_{n-1}) \in \mathbb{F}_2^n \mid c(\gamma^j) = 0 \quad \forall j=1, \dots, d-1 \right\}$$

$c(X) = \sum_{i=0}^{n-1} c_i \cdot X^i$

Notice that this is exactly the same as our def. of RS codes, except that we restrict $(c_0, \dots, c_{n-1}) \in \mathbb{F}_2^n$ instead of \mathbb{F}_2^m . In particular, $\text{dist}(BCH(n, d)) = d$.

b/c $BCH(n, d) \subseteq RS(n, n-d+1)$
 and $RS(n, n-d+1)$ has dist d .

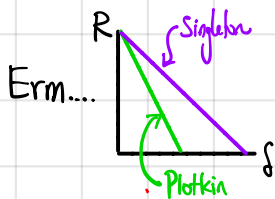
NOTE. BCH codes make sense if you replace "2" with "q" [any prime power].
 We focus on binary codes for today.

FALSE

CLAIM. BCH codes are linear codes with dimension $\geq n - d + 1$.

BAD Proof. $c(\gamma^j) = 0$ is a linear constraint; there are $d-1$ such constraints, so the dimension is at least $n - (d-1)$.

GREAT! So, BCH codes are binary codes that meet the Singleton bound!



doesn't that violate the Plotkin bound??

(Yes).

What's wrong?!

In fact, we can do even better:

CLAIM. BCH(n,d) is \mathbb{F}_2 -linear with dimension $\geq n - \lfloor \frac{d-1}{2} \rfloor \lg(n+1)$.

Proof. We'll show that the linear constraints $c(\gamma^j) = 0$ are actually redundant:
 $c(\gamma^j) = 0 \iff c(\gamma^{2j}) = 0$. This cuts the number of constraints in half, which gives the bound.

SUB CLAIM. For any $c \in \mathbb{F}_2[X]$, $\alpha \in \mathbb{F}_{2^m}$, $c(\alpha) = 0 \iff c(\alpha^2) = 0$.

pf.

$$\begin{aligned} & c(\alpha) = 0 \\ \iff & [e(\alpha)]^2 = 0 \\ \iff & \left[\sum_{i=0}^{n-1} c_i \alpha^i \right]^2 = 0 \\ \iff & \sum_{i=0}^{n-1} c_i^2 \alpha^{2i} = 0 \\ \iff & \sum_{i=0}^{n-1} c_i \alpha^{2i} = 0 \\ \iff & c(\alpha^2) = 0 \end{aligned}$$

Since $0^2 = 0$

Def. of $c(X)$

Since in \mathbb{F}_{2^m} , $1+1=0$, so $(a+b)^2 = a^2 + ab + ab + b^2 = a^2 + b^2$

Since $c_i \in \{0,1\}$, so $c_i^2 = c_i$.

Def. of $c(X)$ again.

and more generally
 $(\sum_i a_i)^2 = \sum_i a_i^2$.

and by the above reasoning the **SUB CLAIM** proves the **CLAIM**.

Is this good? $R \geq \frac{n - \lfloor \frac{d-1}{2} \rfloor \cdot \lg(n+1)}{n} \sim 1 - \frac{\delta}{2} \lg(n)$, aka

BCH codes

Rate R , distance $\delta \sim \frac{2}{\lg(n)} \cdot (1-R)$.

So BCH codes do slightly better than our STRAWMAN:

STRAWMAN

Rate R , distance $\delta \sim \frac{1-R}{\lg(n)}$

But still not asymptotically good ".

However! Note that BCH codes can be decoded with either Berlekamp-Welch or Berlekamp-Massey!
 So that's pretty cool.

② BINARY REED-MULLER CODES

(Dumb) idea: just do RS codes over \mathbb{F}_2 directly! $RS_2(n, k) = \left\{ (f(\alpha_1), \dots, f(\alpha_n)) \mid \begin{matrix} f \in \mathbb{F}_2[X] \\ \deg(f) < k \end{matrix} \right\}$

This is obviously dumb since (a) $\deg(f) \leq q-1 = 1$ to be interesting

(b) $\alpha_1, \dots, \alpha_n$ should be distinct pts in \mathbb{F}_2 , so $n \leq 2$.

However, one fix is to add more variables.

DEF. The BINARY m -VARIATE REED-MULLER CODE of DEGREE r is

$$RM_2(m, r) = \left\{ \underbrace{(f(\alpha_1), f(\alpha_2), \dots, f(\alpha_{2^m}))}_{\substack{\{\alpha_1, \dots, \alpha_{2^m}\} = \mathbb{F}_2^m, \\ \text{in any pre-determined order.}}} : \underbrace{f \in \mathbb{F}_2[X_1, \dots, X_m]}_{\substack{\text{m-variate polynomials} \\ \text{over } \mathbb{F}_2. \text{ eg,} \\ f(X_1, X_2) = 1 + X_1 X_2 + X_1}}, \underbrace{\deg(f) \leq r}_{\substack{\text{deg means} \\ \text{total degree. eg,} \\ \deg(X_1, X_2) = 2.}}$$

Parameters:

Block length $n = 2^m$

Dimension $k = \sum_{j=0}^r \binom{m}{j} = \text{Vol}_2(r, m)$.

Distance $d = ?$

← This is the number of coefficients in

$$f(X_1, \dots, X_m) = \sum_{\substack{S \subseteq [m] \\ |S| \leq r}} c_S \cdot \prod_{i \in S} X_i,$$

a generic degree $\leq r$ m -variate polynomial over \mathbb{F}_2 .

LEMMA (Binary Schwartz-Zippel)

Let $f \in \mathbb{F}_2[X_1, \dots, X_m] \neq 0$, with $\deg(f) \leq r$.

$$\text{Then } \sum_{\alpha \in \mathbb{F}_2^m} \mathbb{1}\{f(\alpha) \neq 0\} \geq 2^{m-r}.$$

We may do the proof later for a more general version, but if you haven't seen this before it's a FUN EXERCISE!

So $\text{dist}(\text{RM}_2(m,r)) \geq 2^{m-r}$. This is because $\text{RM}_2(m,r)$ is linear, and so as usual we only need to look at the minimum wt of a codeword.

And, it turns out this is the correct answer: consider $f(x_1, \dots, x_m) = x_1 \cdot x_2 \cdot \dots \cdot x_r$. This vanishes whenever any of $x_1, \dots, x_r = 0$, and so

$$|\{ \alpha \in \mathbb{F}_2^m : f(\alpha) \neq 0 \}| = |\{ \alpha \in \mathbb{F}_2^m : \alpha_1 = \dots = \alpha_r = 1 \}| = 2^{m-r}.$$

So for $\text{RM}_2(m,r)$:

$$\begin{aligned} n &= 2^m \\ k &= \text{Vol}_2(r, m) & \Rightarrow & R = \text{Vol}_2(r, m) / 2^m \\ d &= 2^{m-r} & & S = 1/2^r \end{aligned}$$

RM codes also admit efficient decoding algs. We'll see some of these later in the course.

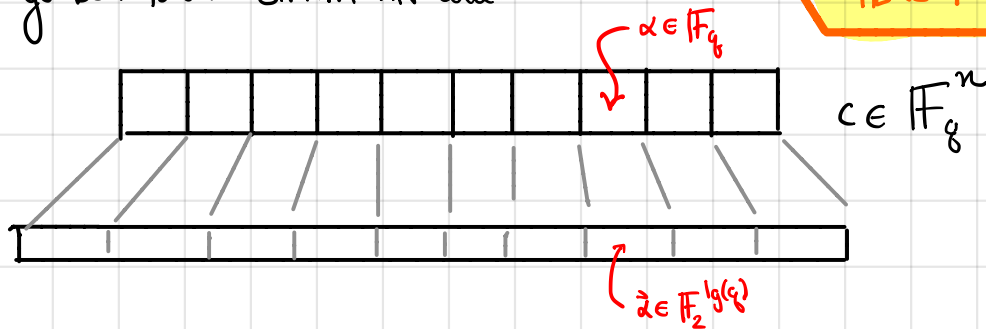
Unfortunately, this isn't asymptotically good either. If $S = \Theta(1)$, then r is constant but $m \uparrow$, so $R \downarrow 0$.

So this doesn't achieve the GOAL either... \cap

③ CONCATENATED CODES.

NOTE: WE BARELY STARTED TALKING ABOUT THESE IN CLASS, AND WILL PICK UP HERE NEXT TIME.

Let's go back to our STRAWMAN code:

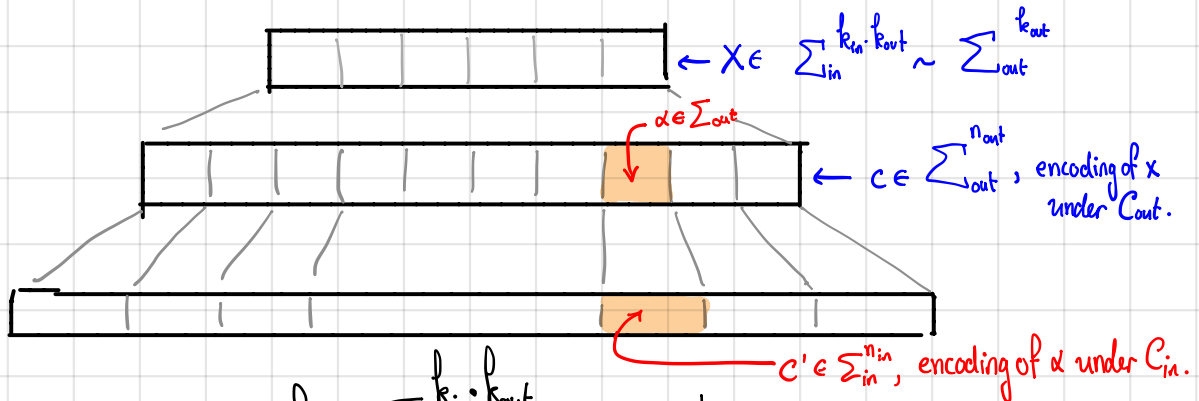


This wasn't a good idea, because if I were a bad guy, I'd mess up one bit from each of $d+1$ of the length $lg(q)$ blocks. So these $\boxed{}$ blocks were not very robust.

IDEA. Let's encode each block $zeta \in \mathbb{F}_2^{lg(q)}$ with a good binary ECC!

DEF. Let $C_{out} \subseteq \sum_{out}^{n_{out}}$ be a q_{out} -ary code of dimension k_{out} and distance d_{out} .
Let $C_{in} \subseteq \sum_{in}^{n_{in}}$ be the same thing with "in" subscripts, so that $q_{out} = q_{in}^{k_{in}}$

The CONCATENATED CODE $C_{in} \circ C_{out} \subseteq \sum_{in}^{n_{out} \cdot n_{in}}$ is defined [by picture...] by



Formally, the encoding of $x \in \mathbb{F}_{q_{in}}^{k_{in} \cdot k_{out}}$ is given by:

- Treat $x \in [\mathbb{F}_{q_{in}}^{k_{in}}]^{k_{out}} \cong \mathbb{F}_{q_{out}}^{k_{out}}$

- Let $y = C_{out}(x) \in \mathbb{F}_{q_{out}}^{n_{out}} \cong [\mathbb{F}_{q_{in}}^{k_{in}}]^{n_{out}}$

- Let $c = \underbrace{Enc_{in}(y_1) \circ \dots \circ Enc_{in}(y_{n_{out}})}_{n_{in} \cdot n_{out} \text{ bits}}$ where $Enc_{in}: \mathbb{F}_{q_{in}}^{k_{in}} \rightarrow \mathbb{F}_{q_{in}}^{n_{in}}$, and \circ is concatenation.

Parameters of Concatenated Codes:

alphabet: Σ_{in}

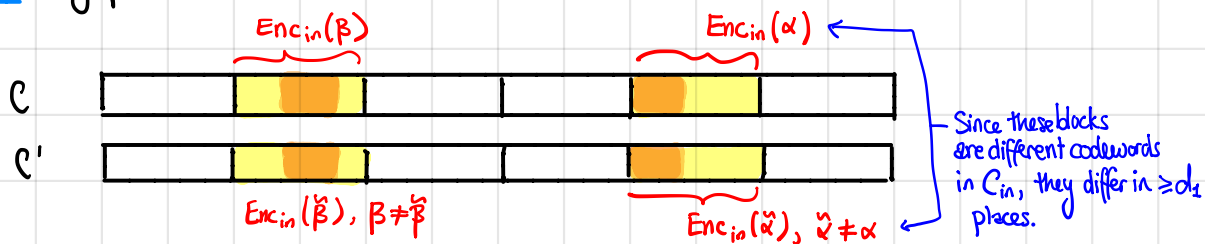
message length: $k_{in} \cdot k_{out}$

codeword length: $n_{out} \cdot n_{in}$

so the rate is $\frac{k_{in} \cdot k_{out}}{n_{in} \cdot n_{out}} = R_{in} \cdot R_{out}$

PROPOSITION. The distance of $C_{in} \circ C_{out}$ is at least $d_{in} \cdot d_{out}$.

pf. by picture: Let $c, c' \in C_{in} \circ C_{out}$:



- At least d_{out} blocks of c, c' are encoding different symbols.
- In each of those, there are at least d_1 symbols in Σ_{in} that differ.
- So that's $d_{out} \cdot d_{in}$ differences total.

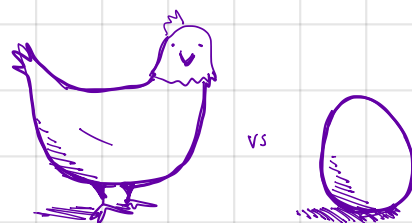
Finally! Progress to our GOAL.

To obtain an EXPLICIT, ASYMPTOTICALLY GOOD BINARY CODES:

1. Set $C_{out} = \text{Reed-Solomon Code}$
2. Set $C_{in} = \text{EXPLICIT, ASYMPTOTICALLY GOOD BINARY CODE}$.

D'oh.

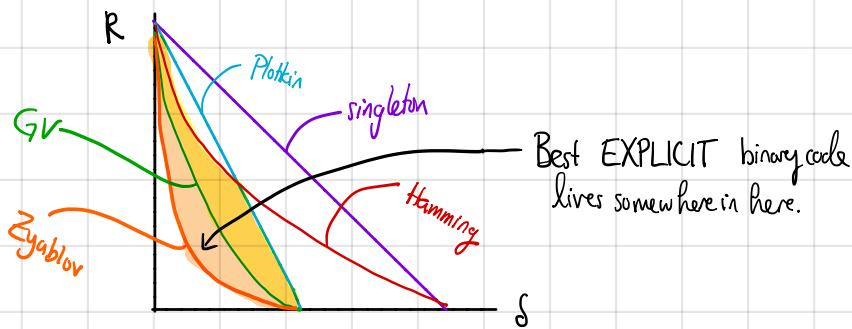
But actually it's OKAY!
The secret is that C_{in} will be short enough that we can do exhaustive stuff efficiently.



THM. For any $\epsilon > 0$, there is a family of explicit binary linear codes of rate R and distance δ , satisfying

$$R \geq \sup_{0 < r < 1 - H_2(\delta) - \epsilon} \left(r \cdot \left(1 - \frac{\delta}{H_2^{-1}(1-r-\epsilon)} \right) \right)$$

That is called the Zyablov Bound.



pf. The construction will be:

- $C_{out} = RS$ code with rate R_{out} , dist. $\delta_{out} = 1 - R_{out}$
- $C_{in} = \underline{\text{Binary linear code on the GV bound}}$, with rate $r \geq 1 - H_2(\delta_{in}) - \epsilon$.

• The concatenated code has:

↖ We still need to say how to get this... we'll get there.

RATE: $R_{out} \cdot r$

DISTANCE: $\delta_{out} \cdot \delta_{in} = (1 - R_{out}) H_2^{-1}(1 - r - \epsilon)$.

Hence $R = R_{out} \cdot r = r \cdot \left(1 - \frac{\delta_{in} \cdot \delta_{out}}{H_2^{-1}(1 - r - \epsilon)} \right)$,

$R_{out} = 1 - \frac{\delta_{in} \delta_{out}}{H_2^{-1}(1 - r - \epsilon)}$

and then we get to choose r .

So that gives us the combinatorial bound.

Next, the algorithmic bit.

continued...

proof continued...

Suppose the evaluation pts for the RS code are \mathbb{F}_q^* , where $q = 2^{k_{in}}$.

So $k_{in} = \lg(q)$, and $n_{out} = q - 1$

ALG 1. Search over all \mathbb{F}_2 -linear codes of rate r and dimension k_{in} .

There are approximately $2^{n_{in} \cdot k_{in}} = 2^{k_{in}^2 / r} = 2^{\lg^2(q) / r}$ such codes

$$q = n_{out} + 1 = \frac{n}{n_{in}} + 1 = \frac{r n}{k_{in}} + 1 \Rightarrow n \sim \frac{1}{r} q \lg(q)$$

So $\lg^2(q) = \Theta(\lg^2(n))$, so that's $2^{\Theta(\lg^2(n))} = n^{\Theta(\lg(n))}$, which is NOT polynomial time. "

ALG 2. On your homework, you will give an alg to construct binary linear codes on the GV bound w/ rate r , dim k_{in} in time $2^{O(k_{in})}$, instead of $2^{O(k_{in}^2)}$. This will fix the above, and proves the theorem.

So we have achieved (most of) our goal.

HOWEVER, this version of "explicit" [can compute it in polynomial time] may be unsatisfying.

WHAT IF I WANTED "explicit" meaning: "Give me a short, useful description" formally, I'd like to be able to compute any entry G_{ij} in time $\text{poly}(\log(n))$.

IDEA: Instead of using the same inner code at every position and requiring it to be good, we'll use a different inner code in each position.

We won't actually know which of these inner codes is good, but we'll know that enough of them are good.

THM. Let $\varepsilon > 0$, fix any k . There is an ensemble of binary linear codes

$$C_{in}^1, C_{in}^2, \dots, C_{in}^N \subseteq \mathbb{F}_2^{2k}$$

of rate $1/2$, with $N = 2^k - 1$, so that for at least $(1-\varepsilon)N$ values of i , C_{in}^i has distance at least $H_2^{-1}(1/2 - \varepsilon)$.

This is called the WOZENCRAFT ENSEMBLE.

proof idea. For $x \in \mathbb{F}_2^k$, treat it as an element of \mathbb{F}_{2^k} .
Then for each $\alpha \in \mathbb{F}_{2^k}^*$, let the α^{th} code C_{in}^α be the image of the encoding map

$$E_{in}^\alpha : x \mapsto (x, \alpha \cdot x)$$

multiplication in \mathbb{F}_{2^k}

treat these as $2k$ bits.

FUN EXERCISE: finish the proof!

Using the Wozencraft ensemble, we can implement the idea above to obtain the JUSTESAN CODE.

DEF (JUSTESEN CODE)

Let $k > 0$ [k will be the dimension of the inner codes in the Wozencraft Ensemble]

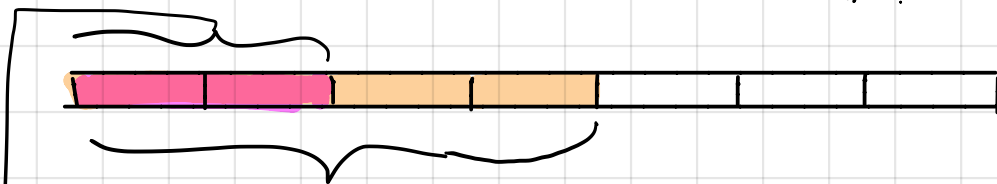
$$\text{Let } C_{\text{out}} = \text{RS}_{2^k}(\mathbb{F}_{2^k}^*, 2^k - 1, R_{\text{out}} \cdot (2^k - 1))$$

Use the Wozencraft Ensemble as the inner code:

$$C = \left\{ \left\langle E_{\text{in}}^\alpha (f(x)) \right\rangle_{\alpha \in \mathbb{F}_{2^k}^*} : f \in \mathbb{F}_{2^k}[X], \deg(f) < R_{\text{out}}(2^k - 1) \right\}$$

CLAIM. Let R_{out} be constant, choose $\varepsilon > \frac{1 - R_{\text{out}}}{2}$. Then \mathcal{C} is asymptotically good.

pf. (sketch) The rate is $R_{\text{out}}/2$, and it's a binary linear code, so we just have to consider the minimum wt to compute the distance. Consider any codeword:



- At least $(1 - R_{\text{out}}) \geq 2\varepsilon$ fraction of the chunks are the encodings of nonzero symbols.

→ • At most an ε -fraction of chunks have "bad" inner codes, so at least an $2\varepsilon - \varepsilon = \varepsilon$ -fraction of chunks are the encodings of nonzero symbols with a "good" inner code.

For each of those, since the inner code has distance $\geq H_2^{-1}(\frac{1}{2} - \varepsilon) = \Theta(1)$, a constant fraction of the bits in each of a constant fraction of blocks are nonzero.

⇒ Each nonzero codeword has relative weight larger than some constant. Thus the code is asymptotically good.

So the JUSTESSEN CODE is "EXPLICIT" in the way we wanted.

The α^i 'th block is given by $(f(x), \alpha \cdot f(x)) \in \mathbb{F}_q^2 \sim \mathbb{F}_2^{2 \lg(q)}$.

That's pretty explicit!

FUN EXERCISE. What is the best rate/distance trade-off you can get w/ the Justesen code?

QUESTIONS to PONDER

- ① How would you efficiently decode a concatenated code?
- ② How would you efficiently decode Reed-Muller codes.
- ③ What happens to the Wozencraft ensemble if you do

$$x \mapsto (x, \alpha \cdot x, \alpha^2 \cdot x, \dots, \alpha^r \cdot x) ?$$