# CS250/EE387 – LECTURE 8 – APPLICATIONS of CONCATENATED RS CODES

## AGENDA

① Syndrome decoding, sparse recovery, group testing

② Application: group testing

① The SYNDROME DECODING problem that we've seen a few times now is:

$n-k$ { [ parity-check matrix $H$ ] (width $n$) } × [ sparse vector $e \in \mathbb{F}_q^n$ ] = [ ]

← $H \cdot e = H(c+e)$ is the SYNDROME.

(over $\mathbb{F}_q$)

PROBLEM: Given $H \cdot e$, recover $e$.

GOAL: Make $n-k$ as small as possible.

This set-up might look familiar.

**ONE SYNTACTICALLY SIMILAR PROBLEM:** SPARSE RECOVERY / COMPRESSED SENSING.

$m$ { [ sensing matrix $\underline{\Phi} \in \mathbb{R}^{m \times n}$ ] (width $n$) } × [ sparse vector $x \in \mathbb{R}^n$ ] = [ ]

← $\Phi x = y$ = "OBSERVATIONS" $\in \mathbb{R}^m$

PROBLEM: Given $\Phi x$, recover $x$.

GOAL: Make $m$ as small as possible.

# Why might we care about this?

1. Image processing and signal processing.



Image (not sparse)    →    Appropriate change-of-basis

· Most natural images/signals are sparse(ish) in some basis  (or w/rt some dictionary).
· So if we can acquire that image/signal by just measuring linear combinations of it and storing those, we can save time and space.

2. Streaming algorithms:

Consider a data stream:

$$X_1, X_2, X_3, \ldots, X_t, \ldots \in \text{ some universe } \mathcal{U} \text{ of size } n$$

You are interested in the frequency counts $f_i$ = #times $i \in \mathcal{U}$ showed up.

But you don't want to store the vector $f \in \mathbb{R}^n$, especially if only a few items show up often.

Instead, keep a SKETCH



= sketch.

f is sparse (ish)

When a new item arrives, you can update the sketch by adding the appropriate column of $\Phi$.

So this is exactly the same as syndrome decoding, except over $\mathbb{R}$ instead of $\mathbb{F}$.

# ANOTHER SYNTACTICALLY SIMILAR PROBLEM: GROUP TESTING.

Let $\mathcal{B} = \{0,1\}$, with the operations "$+$" = $\vee$ (aka, OR) and "$*$" = $\wedge$ (aka AND).

$$m \left\{ \begin{array}{|c|} \hline \text{Pooling matrix} \\ \underline{\Phi} \in \mathcal{B}^n \\ \hline \end{array} \right. \quad \boxed{\phantom{x}} = \boxed{\phantom{x}} \quad \leftarrow \underline{\Phi} \cdot x \in \mathcal{B}^m = \text{"test outcomes."}$$

$\leftarrow$ sparse vector $x \in \mathcal{B}^n$

**PROBLEM:** Given $\underline{\Phi} \cdot x$, recover $x$.

**GOAL:** make $m$ as small as possible.

## Why might we care about this?

Suppose there are $n$ pots of coffee:

- Unfortunately, $s \ll n$ of them are poisoned, but we don't know which.
- Fortunately, there are many lab rats available.[*] If a lab rat has even a drop of poisoned coffee today, then tomorrow they will be sick.[**]

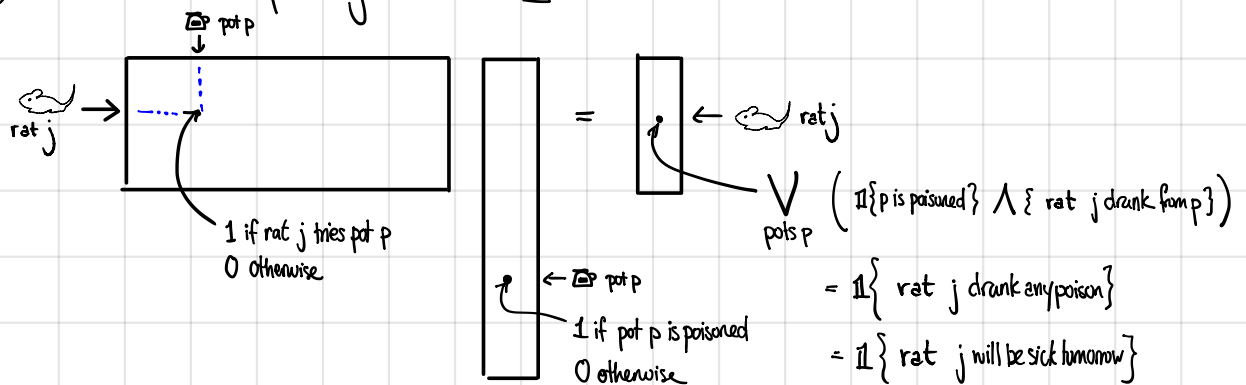  [*] That is, "borrowed" from the biology department...

- You want to decide BY TOMORROW which pots of coffee are poisoned (so that you can drink the rest), while using as few lab rats as possible (so that the biology dept. doesn't notice...)
- The idea is to POOL the samples of coffee:

$n$ pots:

$m$ rats:

- If a lab rat drinks from ANY poisoned coffee pot, they become sick.[**]

[**] Not THAT sick. No animals were harmed in the making of these lecture notes.

Thus, if we make a "pooling matrix" $\Phi$ as below, we have (over $\mathbb{B}$):



↓ pot p

rat j →

1 if rat j tries pot p
0 otherwise

=

← rat j

← pot p

1 if pot p is poisoned
0 otherwise

$\bigvee_{\text{pots } p} \left( \mathbb{1}\{p \text{ is poisoned}\} \wedge \{ \text{rat } j \text{ drank from } p \} \right)$

$= \mathbb{1}\{ \text{rat } j \text{ drank any poison} \}$

$= \mathbb{1}\{ \text{rat } j \text{ will be sick tomorrow} \}$

So that's the picture we had before.

The PROBLEM is to recover $x$, the indicator vector of poisoned pots,
and the GOAL is to minimize the number of lab rats "borrowed" from the biologists.

MORE SERIOUSLY, this problem is usually motivated as follows:

- During WWII, the problem was introduced for testing US soldiers for syphilis.

soldiers ⟷ coffee pots
blood sample ⟷ coffee sample
syphilis tests ⟷ lab rats

- Nowadays, for high-throughput screening.

civilians ⟷ coffeepots
DNA samples ⟷ coffee sample
genetic tests ⟷ lab rats

Tests are expensive, and not many soldiers/civilians are sick, so we'd like to use as few tests as possible.

So both GROUP TESTING and COMPRESSED SENSING are syntactically very similar to SYNDROME DECODING, it's just that they happen over $\mathbb{B}$, $\mathbb{R}$ or $\mathbb{C}$, and $\mathbb{F}_q$, respectively.
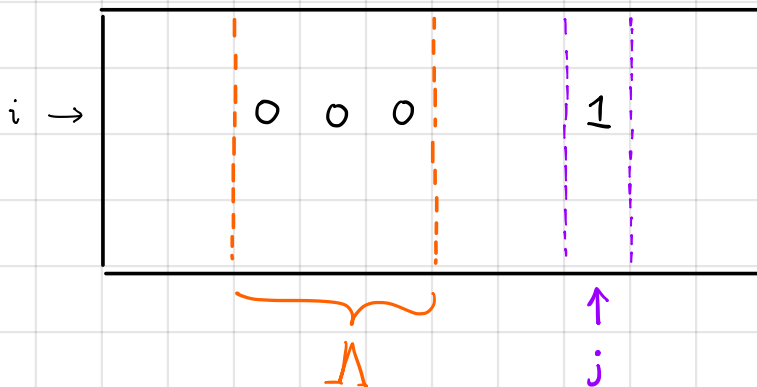
The different algebraic and geometric structures make these problems very different. However, ideas from one are often useful in others.

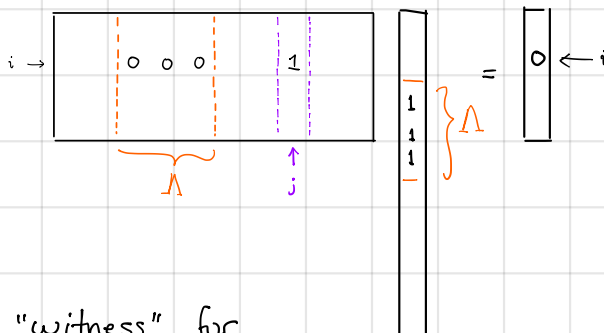Today, we'll see how RS codes can be used to make good GROUP TESTING matrices.

<span style="color:red">↳ NOTE THE CHANGE to N. This will avoid notational collisions later.</span>

**DEF.** A pooling matrix $\underline{\Phi} \in \mathbb{B}^{m \times N}$ is $d$-disjunct if for any set $\Lambda \subseteq [N]$ of size $d$, and any $j \in [N] \setminus \Lambda$, there is at least one $i \in [m]$ so that :

$$\underline{\Phi}_{ij} = 1 \quad \text{and} \quad \overline{\underline{\Phi}}_{i\ell} = 0 \quad \forall\, \ell \in \Lambda$$

Picture:



This is a good thing b/c if $\Lambda$ we the true set of positives (aka, poisoned coffeepots), Then



which gives a "witness" for j's status as not-poisoned.

**THM.** If $\underline{\Phi}$ is d-disjunct, then as a pooling design it can identify up to d positive items.

Moreover, there is an algorithm that runs in time $O(m \cdot N)$ to identify the d positives.

**Pf.** The algorithm is:

    for each j ∈ [N]:
        if all the tests that j participates in are positive:
            label j as positive.
        else  j is not positive.

Why does this work?   Suppose that $\Lambda$ is the true positive set and $j \notin \Lambda$. Then the def of d-disjunctness says that some test $i \in [m]$ which j participates in will come up negative, so the alg will label j "**NOT POSITIVE**." OTOH, if $j \in \Lambda$, then by def. every test it participates in will be positive, so the alg will label j "**POSITIVE**."

So the goal is to come up with d-disjunct matrices $\Phi \in \mathbb{B}^{m \times N}$ so that m is as small as possible.
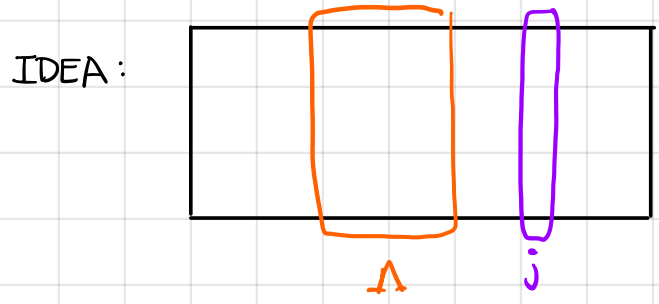
**BEST CONSTRUCTIONS KNOWN:**    $m = O\left(d^2 \log^2_d(N)\right)$    [Kautz-Singleton '64] – we'll see this today
(based on RS codes)

$m = O\left(d^2 \log(N)\right)$    A random matrix does this – or check out [Porat-Rothschild '08] for an explicit construction. (also based on coding theory).

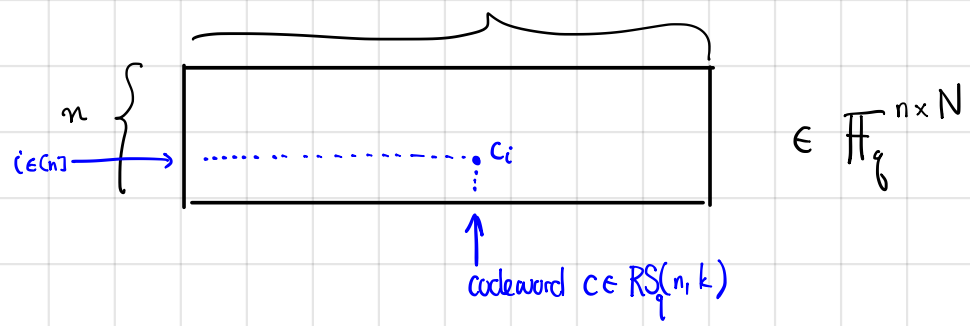BEST LOWER BOUNDS: $m = \Omega\left( d^2 \log_d(N) \right)$ [Dyachkov-Rykov '86]

ALGORITHMS: If $m = O(d^2 \log(N))$, there's an EXPLICIT construction w/ SUBLINEAR TIME algorithm. [Ngo-Porat-Rudra '11(?)] (Also based on coding theory).
We'll see some faster algs later in the course.

Today: A construction with $m = O\left( d^2 \log_d^2(N) \right)$.

IDEA:



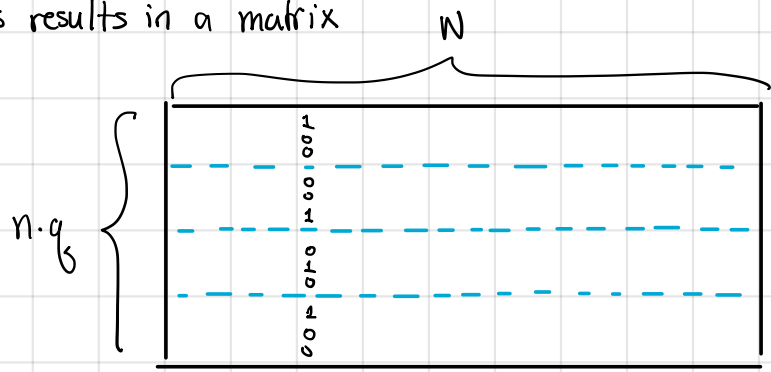We'd like all these columns to be kind of far apart ... let's use codewords!

Let $C = RS_q(n, k)$, let $N = q^k$, $m = q \cdot n$. Consider the matrix formed by:



$\in \mathbb{F}_q^{n \times N}$

codeword $c \in RS_q(n, k)$

Now replace each symbol $\alpha \in \mathbb{F}_q$ w/ a vector of length $q$.

$$\alpha_1 \longleftrightarrow \begin{pmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \quad \alpha_2 \longleftrightarrow \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \quad \ldots, \quad \alpha_q \longleftrightarrow \begin{pmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{pmatrix}$$ where $\mathbb{F}_q = \{\alpha_1, -, \alpha_q\}$.

This results in a matrix



aka, we've CONCATENATED $RS_q(n, k)$ with the code $\alpha_i \in \mathbb{F}_q \longmapsto e_i \in \mathbb{F}_2^q$
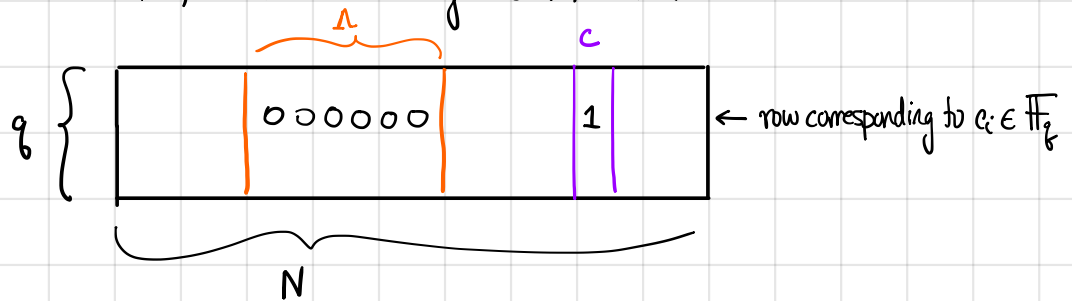
**THM.** If $\mathrm{dist}(\mathcal{C}) > n \cdot \left(\frac{d-1}{d}\right)$, then the matrix $\Phi$ obtained this way is $d$-disjunct.
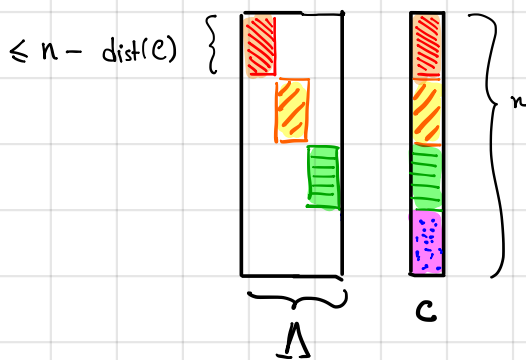
**Proof** (by picture)

Because of how the construction works, we need to show:

$$\forall \; \Lambda \subseteq \mathcal{C}, \; |\Lambda| \leq d, \; \forall \; c \in \mathcal{C} \setminus \Lambda, \; \exists i \in [n] \text{ s.t. } c_i \notin \{w_i : w \in \Lambda\}.$$

Indeed, if that were true, then the $i^{th}$ layer would look like



$\leftarrow$ row corresponding to $c_i \in \mathbb{F}_q$

So consider any $\Lambda \subseteq \mathcal{C}$, $|\Lambda| \leq d$, and any $c \in \mathcal{C} \setminus \Lambda$.



$\leq n - \mathrm{dist}(\mathcal{C})$

· The first column of $\Lambda$ agrees w/ $c$ in at most $n - \mathrm{dist}(\mathcal{C})$ places:  these ones in the picture
· The second column of $\Lambda$ agrees w/ $c$ (and not w/ the first col) in $\leq n - \mathrm{dist}(\mathcal{C})$ places.  these ones
· etc.

Altogether, there are at most $|\Lambda| \cdot (n - \mathrm{dist}(c))$ positions of $c$ that are agreed with by SOME column of $\Lambda$.
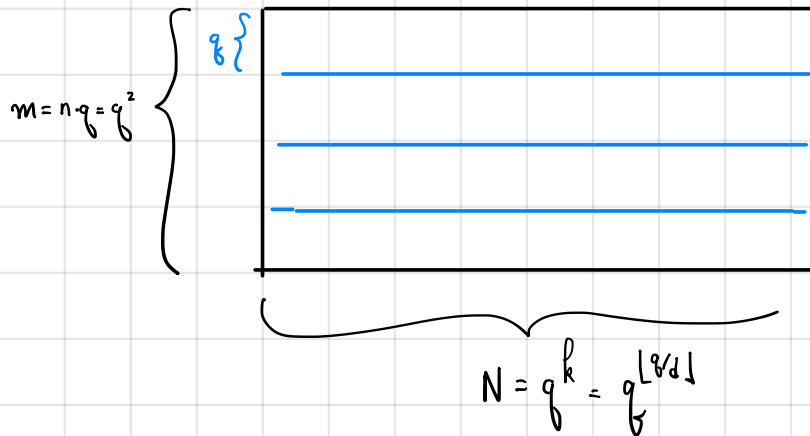
By our guarantee on $\mathrm{dist}(\mathcal{C})$, $\quad |\Lambda|(n - \mathrm{dist}(\mathcal{C})) < d\left(n - n\left(\frac{d-1}{d}\right)\right) = n$.
So there's at least one position that's not agreed with!

Let's instantiate this with $RS_q(\mathbb{F}_q, q, k)$, so $n = q$, $\text{dist}(C) = q - k + 1$

Setting $\text{dist}(C) = n\left(\frac{d-1}{d}\right) + 1 = q\left(\frac{d-1}{d}\right) + 1$, we get $k = \lfloor q/d \rfloor$.

Then our matrix is:



$$m = n \cdot q = q^2$$

$$q\{$$

$$N = q^k = q^{\lfloor q/d \rfloor}$$

Thus we choose $q = \sqrt{m}$, which implies $\log_q(N) = \left\lfloor \frac{\sqrt{m}}{d} \right\rfloor$ aka, $\sqrt{m} \approx d \log_q(N)$.

Then $m \approx d^2 \frac{\log^2(N)}{\frac{1}{2}\log^2(m)}$ which implies $m = O\left(d^2 \frac{\log^2(N)}{\log^2(d)}\right)$,

as claimed.

③ QUICK NOTE ABOUT COMPRESSED SENSING.

A very similar construction can be used to get deterministic compressed sensing matrices.

For those who know the lingo, this EXACT SAME construction $^{\text{appropriately normalized}}$ is an S-RIP matrix $\Phi \in \mathbb{R}^{m \times N}$ with $m = O(s^2 \log^2(N))$.

And you can do slightly better if you replace $\mathbb{F}_p$ w/ the $p^{th}$ roots of unity.

[ See Cheraghchi's "Coding-Theoretic Methods for Sparse Recovery" for lots more! ]

# QUESTIONS TO PONDER

① Can you come up with a recovery scheme for this group testing matrix that runs in time $poly(d\log(N))$ [in particular, sublinear in n?]

② Can you make a group testing $\overset{\smash{\text{↙ or compressed sensing}}}{\text{scheme}}$ using the semantic similarity to syndrome decoding? (Rather than the scheme we saw, which used a different connection to coding theory)