# CS 250/ EE 387 - LECTURE 3 -

## AGENDA

① GV Bound
② Efficiency?
③ Detour: McEliece crypto system
④ Off to asymptopia!
⑤ Q-ary entropy

> ### TODAY'S OCTOPUS FACT
>
> Octopuses have no bones, and can squeeze through very small spaces.
> In 2016, an octopus named "Inky" in a New Zealand aquarium escaped to the ocean through a 6-inch diameter drainpipe.
>
> Freedom!

---

⓪ Recap.    Last time, we saw LINEAR CODES:

$$C = \left\{ \;^\wedge\left[\;\boxed{G}\;\right]\genfrac{}{}{0pt}{}{[]}{x} \; : \; x \in \mathbb{F}_q^k \right\} = \left\{ c \in \mathbb{F}_q^n \; : \; \left[\;\boxed{H}\;\right]\genfrac{}{}{0pt}{}{[]}{c} = 0 \right\}$$

$\underset{\text{G is a GENERATOR MATRIX}}{\Big\uparrow}$ $\qquad$ $\underset{\substack{\text{H is a}\\\text{PARITY-CHECK MATRIX.}}}{\Big\uparrow}$

$k$ (over G)

- Linear algebra works great if $\mathbb{F}_q$ is a field.
- This is super useful.

Today: A few useful things one can do with linear codes:

- GILBERT-VARSHAMOV BOUND
- McELIECE CRYPTOSYSTEM. ← Technically, this is a useful thing you CAN'T do with linear codes.

(Plus, a few combinatorial things you can't do with any codes, useful or otherwise).

# (1) The GILBERT VARSHAMOV BOUND

So far, we have seen the **HAMMING BOUND**, which is an upper bound on the rate of a code. (aka, an IMPOSSIBILITY RESULT). We can match it for $n = 7$, $k = 4$, but what about in general?

Next, we'll see the **GILBERT-VARSHAMOV BOUND**, which is a **POSSIBILITY RESULT**.

---

**THM** (GILBERT-VARSHAMOV)

For any prime power $q$, and for any $d < n$, there exists a linear code $C$ of length $n$, alphabet size $q$, distance $d$, and rate

$$R \geq 1 - \frac{\log_q (\text{Vol}_q (d-1, n)) - 1}{n}$$

NOTE: You can remove the words "prime power" and "linear" and the statement is still true.

---

Compare this to the Hamming bound, which said:

$$R \leq 1 - \frac{\log_q (\text{Vol}_q (\lfloor \frac{d-1}{2} \rfloor, n))}{n}$$

The only difference is this part.

We will talk more about the relationship between these two later, but for now just notice that $R_{GV} < R_{HAMMING}$, so math is not broken.

We'll prove the GV bound now — it's pretty easy! However, it's **NOT KNOWN** if we can do better in general!

**OPEN**

---

**QUESTION**    Do there exist binary codes that do better than the GV bound? (for all parameter regimes?)

Proof of the GV bound.

IDEA : A <u>random</u> linear code will do the trick with probability $> 0$.
So, in particular, there exists a linear code that works!

Let $C$ be a random subspace of $\mathbb{F}_q^n$, of dimension $k$. ⟵ $k$ TBD
Let $G$ be a random generator matrix for $C$.

USEFUL FACT: For any fixed $x \neq 0$, $G \cdot x$ is uniformly random in $\mathbb{F}_q^n \setminus \{0\}$.

Informal proof: Because of all the symmetry, how could it be any other way?
formal proof: Fun exercise!

Now, $\text{dist}(C) = \min_{c \in C \setminus \{0\}} \text{wt}(c) = \min_{x \in \mathbb{F}_q^k \setminus \{0\}} \text{wt}(G \cdot x)$

For any given $x \neq 0$, by the USEFUL FACT,

$$\mathbb{P}_G\left\{ \text{wt}(G \cdot x) < d \right\} = \mathbb{P}_G\left\{ G \cdot x \in B_q(0, d-1) \right\}$$

$$\leq \frac{\text{Vol}_q(d-1, n)}{q^n}$$

⟵ number of choices for $G \cdot x$ w/ $\text{wt} \leq d-1$

⟵ number of choices for $G \cdot x$ total.

(Technically it's $\frac{\text{Vol}(d-1, n) - 1}{q^n - 1}$, which is $\leq$ the above)

So by the union bound,

$$\mathbb{P}\left\{ \exists x \in \mathbb{F}_q^k : \text{wt}(G \cdot x) < d \right\} \leq q^k \cdot \frac{\text{Vol}_q(d-1, n)}{q^n}.$$

Thus, we win as long as this is $\not\leq 1$. Taking logs of both sides, we win if

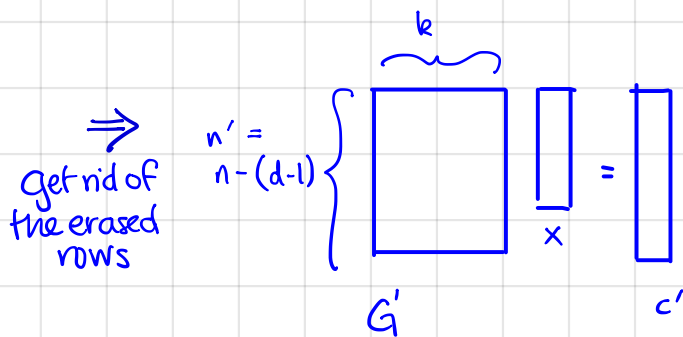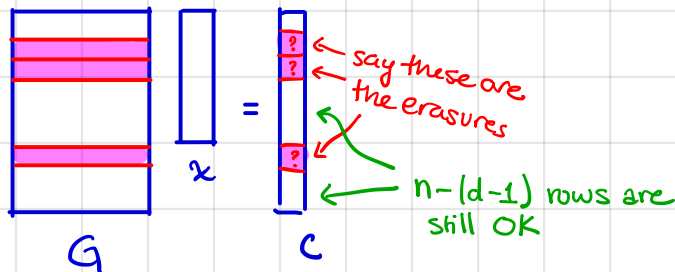$$k - n + \log_q\left( \text{Vol}_q(d-1, n) \right) < 0$$

So choose $k = n - \log_q(\text{Vol}_q(d-1, n)) - 1$, and we are done. ∎

② EFFICIENCY (?)

- If $C$ is linear, we have an efficient encoding map $x \longmapsto G \cdot x$

  The computational cost is one matrix-vector multiply

- If $C$ is linear with distance $d$, we can DETECT $\leq d-1$ errors efficiently:

  If $0 < wt(e) \leq d-1$ and $c \in C$, then $H(c+e) = H \cdot e \neq 0$, so just check if $H\check{c} \neq 0$.

- If $C$ is linear with distance $d$, we can CORRECT $\leq d-1$ ERASURES efficiently:

  We have



  ← say these are the erasures

  ← $n-(d-1)$ rows are still OK

  $G$     $x$     $C$

$\Rightarrow$ Get rid of the erased rows

$n' = n-(d-1)$



  $k$

  $G'$     $x$     $c'$

  Because the distance is $d$, there is exactly one $c' \in \mathbb{F}^{n'}$ that is consistent with these equations, and hence $G'$ is full rank.

$\Rightarrow$ Solve this linear system $G'x = c'$ for $x$.

- If $C$ is linear with distance $d$, can we CORRECT $\lfloor \frac{d-1}{2} \rfloor$ ERRORS efficiently?

- It worked for the $(7,4,3)_2$-Hadamard code!
- But what about in general?

- Consider the following problem:

$$\text{Given } \check{c} \in \mathbb{F}_q^n, \text{ and } G \in \mathbb{F}_q^{n \times k}, \text{ find } x \in \mathbb{F}_q^k \text{ s.t. } \Delta(G \cdot x, \check{c}) \text{ is minimized.}$$

aka, find the codeword closest to a received word $\check{c}$.

- This problem (called MAXIMUM-LIKELIHOOD DECODING for LINEAR CODES) is NP-hard in general [Berlekamp-McEliece-van Tilburg 1978], even if the code is known in advance and you have an arbitrary amount of preprocessing time [Bruck-Noar 1990, Lobstein 1990]. It's even NP-hard to approximate (within a constant factor)! [Arora-Babai-Stern-Sweedyk 1993].

- Even computing the minimum distance of linear codes is NP hard!

- This all sounds bad, but remember that NP-hardness is a worst-case condition. There exist linear codes that are (probably) hard to decode, but that doesn't mean that all of them are.

- We will spend most of this class talking about how to get efficiently-decodable codes. But first let's see one application where the hardness is a good thing.

NOTE: Actually, it's not really clear what "NP hard" or "computational efficiency" means here. (Besides the fact that we did not define them).
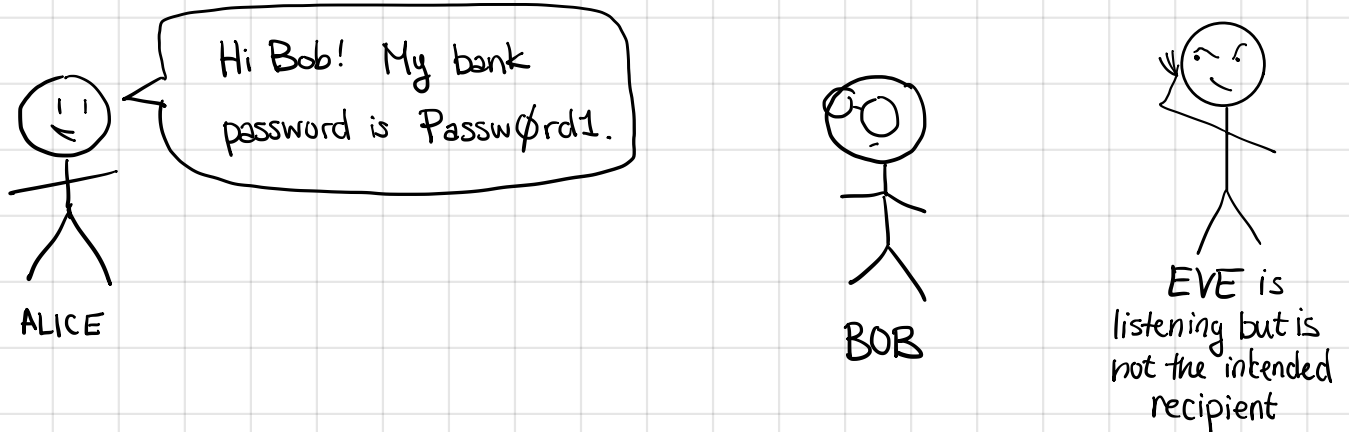In particular, these notions make sense as the input size grows.
What is growing? We'll come back to this in a moment.

\begin {detour}

③ APPLICATION: MCELIECE CRYPTOSYSTEM

Suppose that Alice and Bob want to talk SECURELY.
Now there is no noise, just an EAVESDROPPER Eve.

Hi Bob! My bank
password is Passw∅rd1.

ALICE

BOB

EVE is
listening but is
not the intended
recipient

- In **PUBLIC KEY CRYPTOGRAPHY**, everyone has a public key and a private key.
- To send a message to Bob, Alice encrypts it using Bob's public key.
- Bob decodes it with his private key.
- We hope that this is secure as long as Bob's private key stays private.

**HERE IS SUCH A SCHEME**, using binary linear codes:

- Bob chooses: $G \in \mathbb{F}_2^{n \times k}$ is the generator matrix for an (appropriate) efficiently decodable binary linear code $C$.

  ↑ from $t$ errors

  <span style="color:green">The McEliece system uses something called a "Goppa Code", but we will not go into details in these notes.</span>

- Bob chooses:
  - A random invertible $S \in \mathbb{F}_2^{k \times k}$
  - A random permutation matrix $P \in \mathbb{F}_2^{n \times n}$

- Bob's private key: $(S, G, P)$

- Bob's public key: $\hat{G} = \boxed{\begin{smallmatrix} & & 1 \\ & 1 & \\ & & 1 \\ 1 & & \end{smallmatrix}} \quad \boxed{\phantom{G}} \quad \boxed{\phantom{S}}$ , and $t$

  $\quad\quad\quad\quad\quad\quad\quad\quad\quad\ P \quad\quad G \quad\ S$

To send a message $x \in \mathbb{F}_2^k$ to Bob:
- Alice chooses a random vector $e \in \mathbb{F}_2^n$ with $wt(e) = t$
- Alice sends Bob $\hat{G}x + e$

To decrypt Alice's message $\hat{G}x + e$:
- Bob computes $P^{-1}(\hat{G}x + e) = GSx + P^{-1}e = GSx + e'$, where $wt(e') = t$
- Bob uses his efficient decoder for $C$ to find $S \cdot x$
- Bob computes $x = S^{-1} \cdot Sx$

Why might this be secure?

Suppose Eve sees $\hat{G}x + e$.
She knows $\hat{G}$ and $t$, so this problem is the same as decoding the code $\hat{C} = \{\hat{G}x \mid x \in \mathbb{F}_2^k\}$
from $t$ errors. **WE HOPE THIS IS HARD.**

Note: "Decoding $\hat{G}x + e$ is hard for Eve" is NOT the same as "Maximum likelihood decoding of linear codes is NP-hard."

- First, we have some promise that there were $\leq t$ errors.
- Second, NP-hardness is a worst-case assumption: for crypto we need an average-case assm.

The assm that "Decoding $\hat{G}x + e$ is hard" (for an appropriate choice of $G$) is called the MCELIECE ASSUMPTION. Some people believe it and some don't.

\end{debor}

# ④ OFF TO ASYMPTOPIA

We'll return to computational issues later — but first we need to talk about what we mean by "for large $n$". So for now let's return to the combinatorial question:

## WHAT IS THE BEST TRADE-OFF between RATE and DISTANCE?

So far we've seen two bounds:

$$1 - \tfrac{1}{n} \cdot \log_q \left( \mathrm{Vol}_q (d-1, n) \right) - \tfrac{1}{n} \leq \text{optimal } \tfrac{k}{n} \leq 1 - \tfrac{1}{n} \cdot \log_q \left( \mathrm{Vol}_q (\lfloor \tfrac{d-1}{2} \rfloor, n) \right)$$

Gilbert Varshamov (over left term) — Hamming (over right term) — for $|\Sigma| = q$

So for particular $d, k, n$, these tell us something ... but what do they tell us in general? And what does "in general" mean?

We are going to think about the following limiting parameter regime:

$$\boxed{n, k, d \longrightarrow \infty \quad \text{so that} \quad \tfrac{k}{n} \to R \text{ and } \tfrac{d}{n} \to \delta \text{ approach constants.}}$$

Motivations:
(1) It will allow us to better understand what's possible and what's not
(2) In many applications, $n, k, d$ are pretty large and $R, \delta$ are the things we want to be thinking about.
(3) It will let us talk meaningfully/rigorously about computational complexity.

**DEF** A <u>FAMILY of CODES</u> is a collection $C = \{ C_i \}_{i=1}^{\infty}$, where $C_i$ is an $(n_i, k_i, d_i)_{q_i}$ code.

The <u>RATE</u> of $C$ is $R(C) := \lim\limits_{i \to \infty} k_i/n_i$

The <u>RELATIVE DISTANCE</u> of $C$ is $\delta(C) := \lim\limits_{i \to \infty} d_i/n_i$

**NOTES:**

- We will frequently abuse notation and refer to $C$ as a "code," and we'll drop the subscript $i$ and just think about $n, k, d \to \infty$.

- The alphabet of $C_i$ might depend on $i$. (But if it doesn't will say the whole family has alphabet size $q$).

**EXAMPLE.** In the in-class exercises, you will investigate HAMMING CODES.
The $i^{th}$ code $C_i$ is a $(2^i - 1, 2^i - i - 1, 3)_2$ code.
$C_i$ is defined by its parity-check matrix,



This matrix has every nonzero binary vector as a column.

The RATE of this family is $\lim\limits_{i \to \infty} \dfrac{2^i - i - 1}{2^i - 1} = 1$.

The RELATIVE DISTANCE is $\lim\limits_{i \to \infty} \dfrac{3}{2^i - 1} = 0$.

*Our Question:* What is the best trade-off between $R(C)$ and $\delta(C)$?
(in the asymptotic setting)

*Easier question:* Can we obtain codes with $R(C) \neq 0$ and $\delta(C) \neq 0$?

**DEF.** Such a code (with $R(C) > 0$, $\delta(C) > 0$) is called <u>ASYMPTOTICALLY GOOD</u>.

$$1 - \frac{1}{n} \cdot \log_q\left(\mathrm{Vol}_q(d-1, n)\right) \leq \underset{\text{for } |\Sigma| = q}{\text{optimal } k/n} \leq 1 - \frac{1}{n} \cdot \log_q\left(\mathrm{Vol}_q(\lfloor\frac{d-1}{2}\rfloor, n)\right)$$

Gilbert Varshamov ↓   Hamming ↓

Now that we have an asymptotic parameter regime, how should we parse the GV and Hamming bounds?

In particular, what is $\frac{1}{n}\log_q\left(\mathrm{Vol}_q\left(\lfloor\frac{d-1}{2}\rfloor, n\right)\right)$ in terms of $\delta$, if $\delta \coloneqq d/n$?

OK, so this is $\sum_{j=0}^{\lfloor\frac{d-1}{2}\rfloor} \binom{n}{j}(q-1)^j$, but that is not very helpful.

---

**DEF.**  The q-ary entropy function $H_q : [0,1] \to [0,1]$ is defined by

$$H_q(x) = x\log_q(q-1) - x\log_q(x) - (1-x)\log_q(1-x).$$

---

This generalizes the BINARY ENTROPY FUNCTION $H_2(x) = H(x)$, which you may have seen. The reason we care (for this class) is that $H_q(x)$ captures $\mathrm{Vol}_q(xn)$ nicely:

---

**PROP.**   Let $q \geq 2$ be an integer, and let $0 \leq p \leq 1 - 1/q$.  Then

$$(i) \quad \mathrm{Vol}_q(n, pn) \leq q^{n \cdot H_q(p)}$$
$$(ii) \quad \mathrm{Vol}_q(n, pn) \geq q^{n \cdot H_q(p) - o(n)}$$

A function $f(n)$ is $o(n)$ if $\frac{f(n)}{n} \to 0$ as $n\to\infty$

---
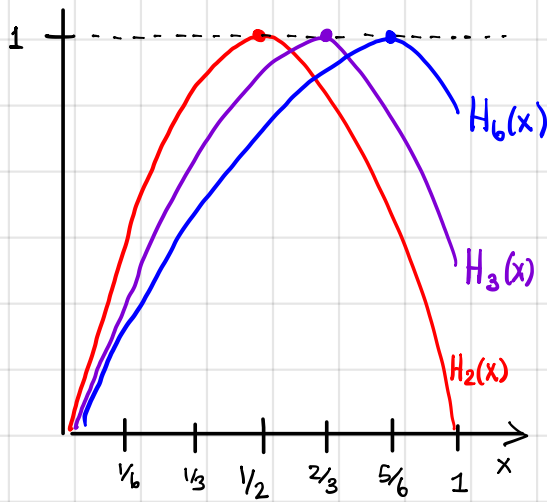
See ESSENTIAL CODING THEORY, Chapter 3 for a proof.
[Proof idea:  mess around with the binomial coefficients and use Stirling's formula.]

$H_q(x)$ looks like this:

↑ something

Some useful properties [You can see these by taking Taylor expansions]:

- If $q$ is really big, (p is reasonable and) then $H_q(p) = p \cdot \log_q(q-1) + \log_q(\text{stuff}) + \log_q(\text{stuff}) \approx p$

  basically 1 ─── really small

  (So, eventually the plot looks like  $\frac{H(x)}{1000000000000}$ )

- If $q$ is reasonable and $p$ is really small, then $H_q(p) = p \cdot \log_q(q-1) + p\log_q(1/p) + (1-p)\log_q(\frac{1}{1-p})$

  $= O(p)$ ─── This term is the largest ─── $\approx P/\ln(q)$ $= O(p)$

  $$\approx p\log_q(1/p)$$

  (So, near $0$ all those curves look like $\frac{x\ln(1/x)}{\ln(q)}$ )

Now, we can take limits in the GV and Hamming bounds to obtain:

**Thm** (HAMMING BOUND) For any family $C$ of $q$-ary codes,
$$R(C) \leq 1 - H_q(\delta(C)/2)$$

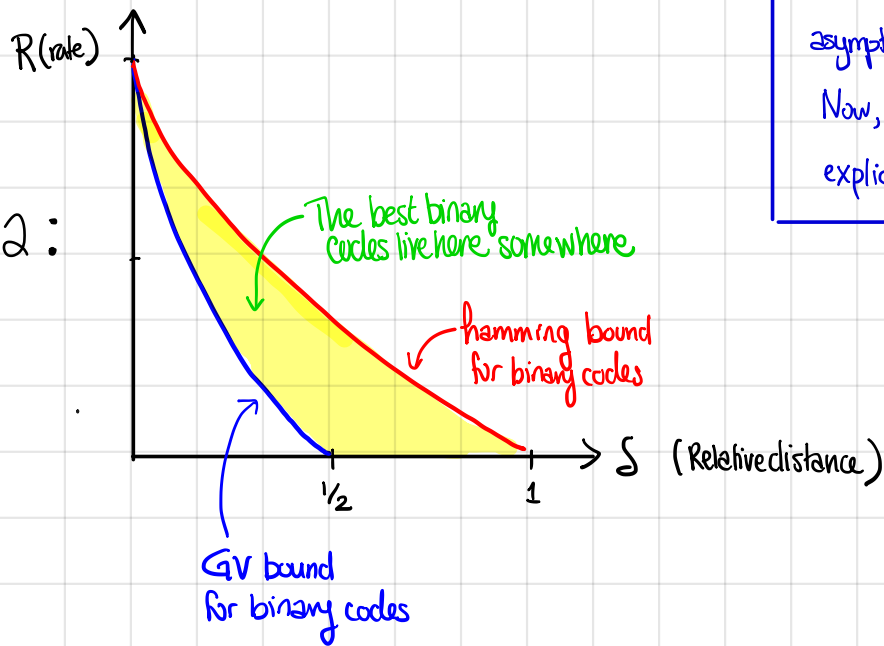Proof: follows by taking limits in the Hamming bound.

**Thm** (GV BOUND) Let $q \geq 2$. For any $0 \leq \delta \leq 1-\frac{1}{q}$, and for any $0 < \varepsilon \leq 1 - H_q(\delta)$, there exists a $q$-ary family of codes $C$ w/ $\delta(C) \geq \delta$ and
$$R(C) \geq 1 - H_q(\delta) - \varepsilon.$$

Proof: FUN EXERCISE (given the nonasymptotic version)

**Now** it's easier to compare these two bounds.



This is the picture for q=2:

- R (rate)
- The best binary codes live here somewhere
- hamming bound for binary codes
- GV bound for binary codes
- $\delta$ (Relative distance)
- ½
- 1

SOME QUESTIONS.

**QUESTION** Are there families of codes that beat the GV bound?

ANSWER 1: Yes. For $q \geq 49$, "Algebraic Geometry Codes" beat the GV bound.

ANSWER 2: ??? For binary codes, we don't know. OPEN PROBLEM!

**QUESTION** Can we find _explicit_ constructions of families of codes that meet the GV bound?

ANSWER 1. For large alphabets, yes. (We'll see soon)

ANSWER 2. ??? For binary codes, recent work of [Ta-Shma 2017] gives something close in a very particular parameter regime... but in general, OPEN PROBLEM!

# QUESTIONS to PONDER

① Can you think of strategies to improve the Hamming bound?

② Is it possible to have a $\overset{\text{family of}}{\vee}$ codes with $\delta > 1 - 1/q$ and $R > 0$?