

CS250/EE387 - LECTURE 1 - LOGISTICS + BASICS

AGENDA

- ① Logistics
- ② Course Pitch
- ③ Basic problem in coding theory
- ④ Formal definitions
- ⑤ Rate vs. Distance: Hamming bound

① LOGISTICS

• COURSE ELEMENTS

- Pre-recorded videos, with corresponding lecture notes
- In-class exercises, meant to practice, reinforce, and extend material in the videos/notes.
- 3 HW assignments
- Final project

• CLASS MEETINGS

- This is a "flipped" class - watch the videos before class and come to class ready to engage!

SEE COURSE WEBSITE FOR MORE DETAILS!

Also for the schedule, materials, assignments, etc.

② COURSE PITCH

"ALGEBRAIC ERROR CORRECTING CODES"

II

I

I. Error correcting codes are a fundamental tool for

II. Algebraic techniques are a fundamental tool for designing ECCs.

- communication
- storage
- complexity theory
- algorithm design
- cryptography
- pseudorandomness
- etc...

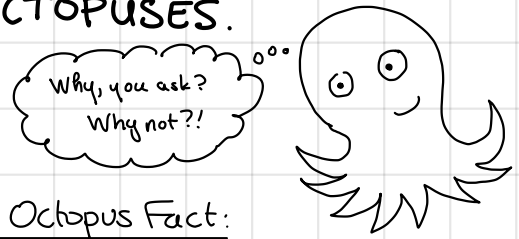
Basically, this course is about the following fact:

LOW-DEGREE POLYNOMIALS

DON'T HAVE TOO MANY ROOTS.

As we will see, this fact is stupidly useful throughout CS and EE.

During this course, ^{at least in the lecture notes...} we will be learning a little bit about **OCTOPUSES**.



Today's Octopus Fact:

Octopuses have three hearts! One of the hearts is inactive when the octopus is swimming, so it tires out faster when swimming than when crawling.

In particular, in these places

In this class we will discuss:

- • Basics of Error Correcting Codes: combinatorial bounds + existential results
- • Some basic abstract algebra [finite fields - nothing fancy]
- • The classic polynomial codes: Reed-Solomon and Reed-Muller
- • [If time we will mention] fancier polynomial codes:
Multiplicity Codes, Folded RS codes
- • Algorithms for manipulating these codes in various settings:
Unique decoding, list decoding, local decoding
- • Applications!

In this class we will **NOT** discuss:

- Nitty gritty details of any one application (this is a THEORY course)
- LDPC codes, Turbo codes, Raptor Codes, Fountain Codes, ...
[See Montanari's course EE 388 for all that good stuff.]

At the end of this course:

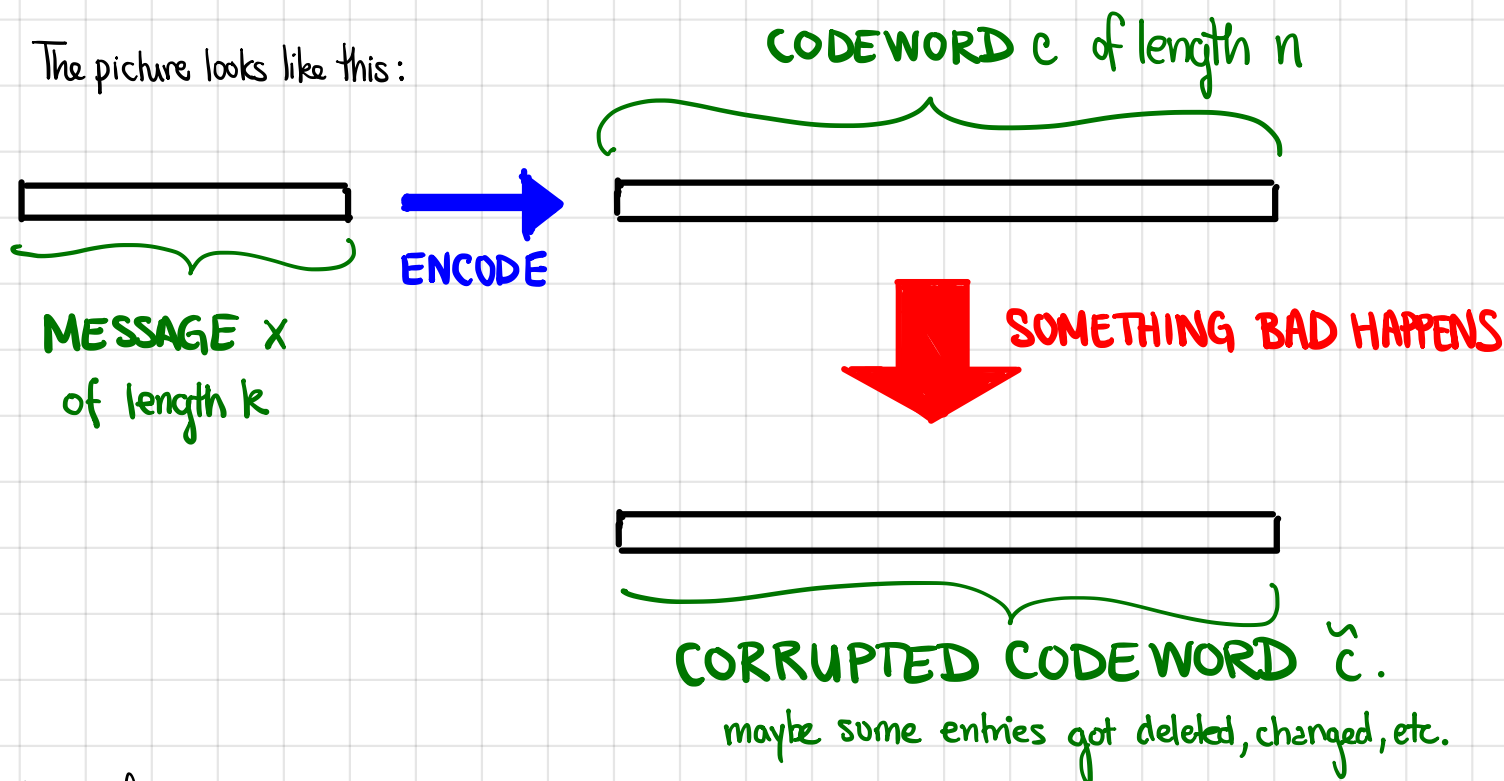
YOU SHOULD HAVE THE TOOLS TO USE
ERROR-CORRECTING CODES (and the algebraic tools behind them)
IN YOUR OWN RESEARCH/LIFE.

That means:

- Enough familiarity with terminology, constructions, algorithms, and notions of decoding to pick up a research paper and understand it.
- Exposure to lots of examples of how ECCs can be useful in a wide variety of settings.

③ The BASIC PROBLEM in CODING THEORY

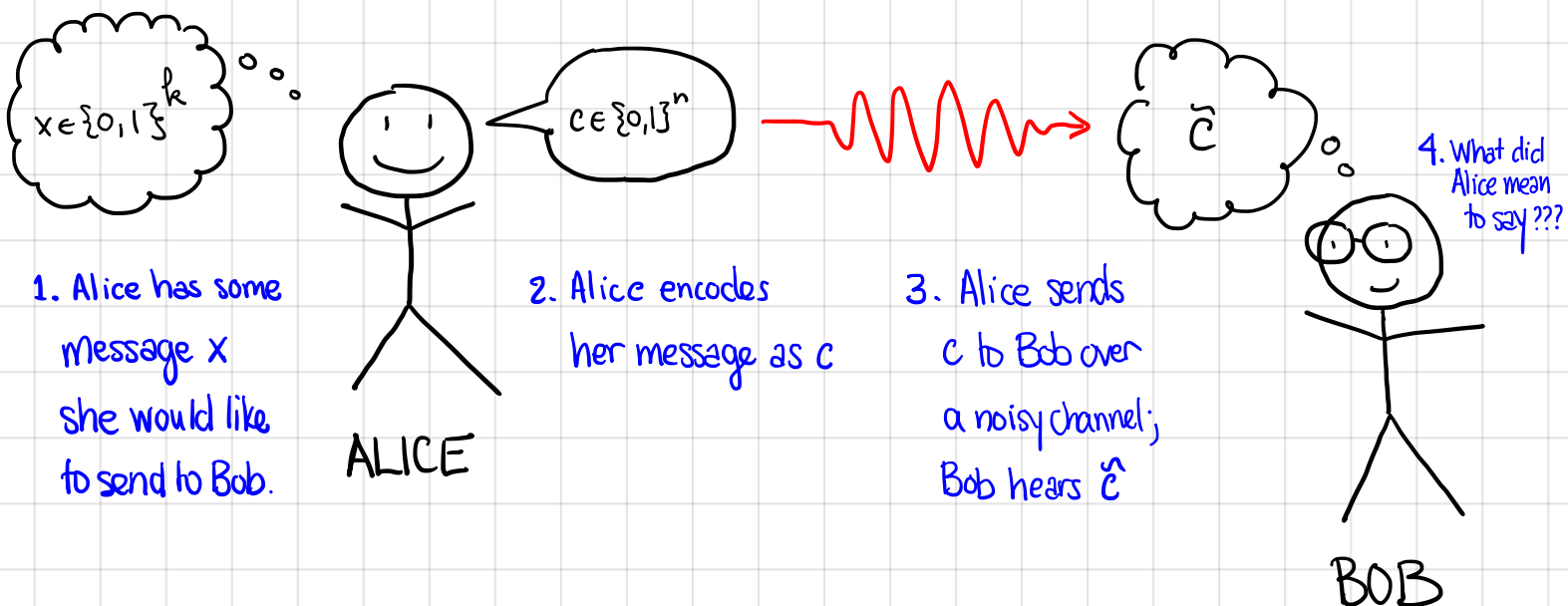
The picture looks like this:



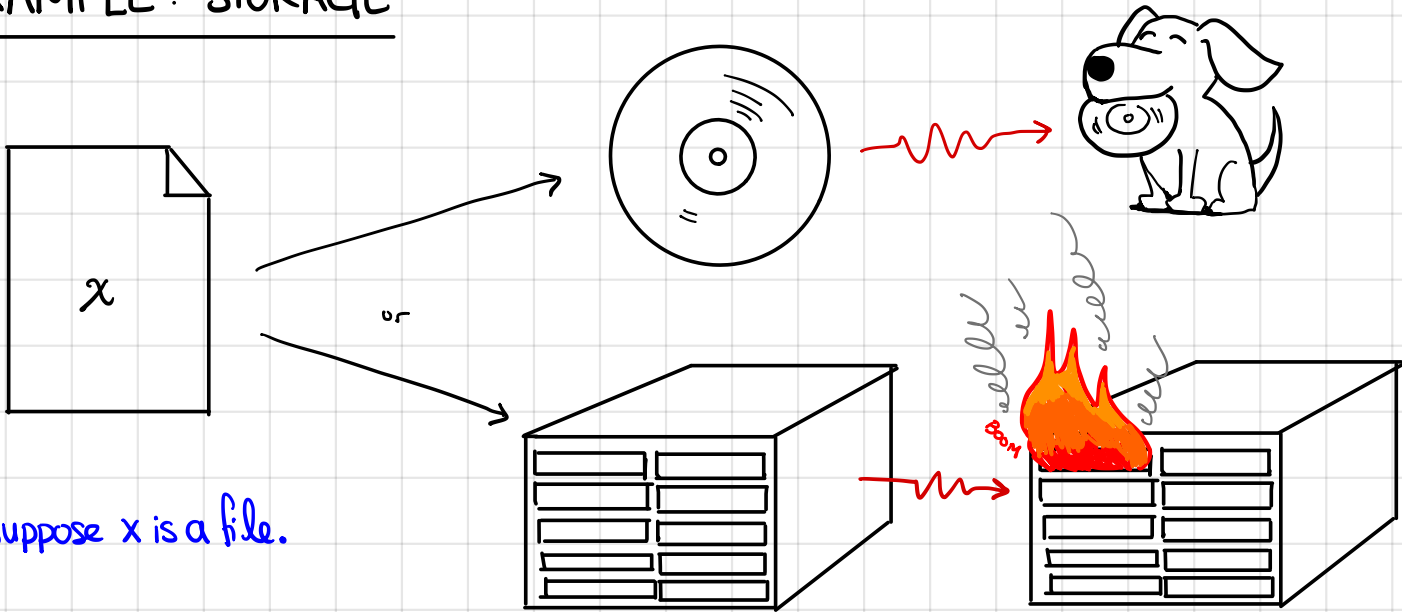
The goal:

GIVEN \hat{c} , FIND (SOMETHING ABOUT) x .

EXAMPLE: COMMUNICATION



EXAMPLE: STORAGE



1. Suppose x is a file.

2. Encode x as a codeword c .

3. c is stored; say on a CD
or in a RAID array. . . .

but something
BAD happens.

5. I still want x !

THINGS WE CARE ABOUT:

- ① We should be able to handle the **SOMETHING BAD**, whatever that means.
- ② We should be able to recover **WHAT WE WANT TO KNOW** about x .
- ③ We want to **MINIMIZE OVERHEAD**: k/n should be as big as possible.
- ④ We want to do all this **EFFICIENTLY**.

QUESTION What are the trade-offs between ①-④?

It depends on how we model things:

- What is the **SOMETHING BAD**?
- What exactly do we **WANT TO KNOW**?
- What counts as **EFFICIENT**? What kind of access do we have to \tilde{c} ?

Today we'll look at one way of answering these questions.

There are many legit ways, and we will see more throughout the quarter.

④ FORMAL DEFINITIONS

Let Σ be any finite set and let $n > 0$ be an integer.

DEF. A CODE C of BLOCK LENGTH n over an ALPHABET Σ is a subset $C \subseteq \Sigma^n$.
An element $c \in C$ is called a CODEWORD.

Sometimes I will say "length" instead of "block length."

So far, this is not a very interesting definition.

EXAMPLE 1. $C = \{ \text{HELLOWORLD, BRUNCHTIME, ALLTHETIME} \}$
is a code of block length 10 over $\Sigma = \{A, B, \dots, X, Y, Z\}$.

EXAMPLE 2.

$$C = \left\{ \begin{array}{l} (0, 0, 0, 0) \\ (0, 0, 1, 1) \\ (0, 1, 0, 1) \\ (0, 1, 1, 0) \\ (1, 0, 0, 1) \\ (1, 0, 1, 0) \\ (1, 1, 0, 0) \\ (1, 1, 1, 1) \end{array} \right\}$$

is a code of length 4 over $\Sigma = \{0, 1\}$.

If $\Sigma = \{0, 1\}$, we say C is a **BINARY CODE**.

This is not a very interesting code, although it does capture the vast majority of my thoughts.

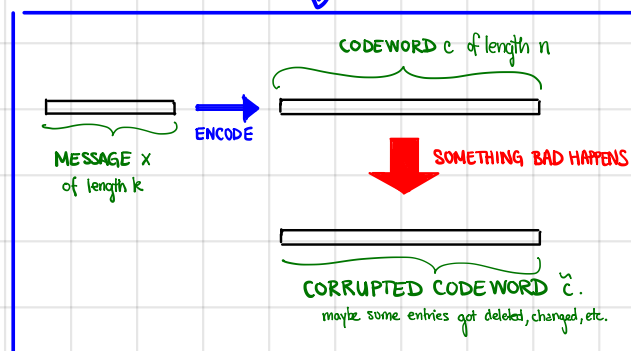
This second example is a bit more interesting.

What does this have to do with the picture from before? [this one]

Consider the map $\text{ENC}: \{0, 1\}^3 \rightarrow \{0, 1\}^4$ given by:

$$\text{ENC}: \underbrace{(x_1, x_2, x_3)}_{\text{message } x} \mapsto \underbrace{(x_1, x_2, x_3, x_1+x_2+x_3 \pmod{2})}_{\text{codeword } c}$$

For example, $\text{ENC}((0, 1, 1)) = (0, 1, 1, 0)$.



Then $C = \text{Im}(\text{ENC})$. That is, C is the set of all codewords that could be obtained using this encoding map.

The second example can actually be used to fix bad stuff.

Suppose you see:

0 ~~X~~ 0 1

← The SOMETHING BAD that happened obscured this entry.

What is the missing bit?

It must be a 1, since $0 + \text{X} + 0 = 1 \pmod{2}$.

← This is called an ERASURE.

We know which bit got erased, but we don't know what its original value was.

Suppose instead you see:

0 0 0 1

Then we know SOMETHING went wrong (at least one bit was flipped) but we are not sure what it was.

← This is called an ERROR.

We say that the code in EXAMPLE 2 can

- ✓ CORRECT one ERASURE
- or -
- ✓ DETECT one ERROR

One bit may have been changed, but we don't know which one.

But it cannot CORRECT one ERROR. Let's see a code that can.

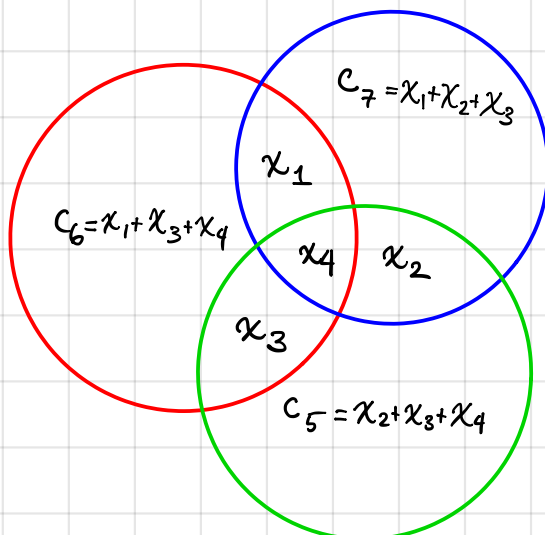
EXAMPLE 3. Consider the encoding map $ENC: \{0,1\}^4 \rightarrow \{0,1\}^7$

↙ all mod 2

$$ENC: (x_1, x_2, x_3, x_4) \mapsto (x_1, x_2, x_3, x_4, x_2+x_3+x_4, x_1+x_3+x_4, x_1+x_2+x_4)$$

Let $C = \text{Im}(ENC)$. So $C \subseteq \{0,1\}^7$, aka C is a BINARY CODE of LENGTH 7.

ANOTHER WAY to VISUALIZE this CODE:



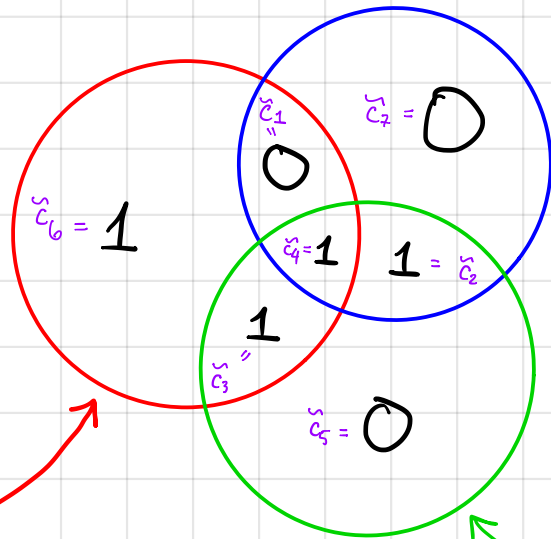
← Put the message x_1, x_2, x_3, x_4 in the middle, and then the circles tell you how to fill in the rest.

PUZZLE: I took some $c \in \mathcal{C}$ and flipped at most one bit, to obtain:

$$\tilde{c} = (0 \ 1 \ 1 \ 1 \ 0 \ 1 \ 0)$$

What is c ?

ALTERNATIVE LOOK
at the PUZZLE:

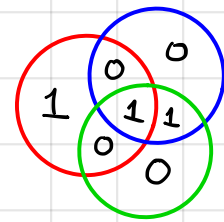


SOLUTION:

In a legit codeword, all three circles should sum to 0. Here, both the green circle and the red circle are messed up.

The thing those circles have in common is \tilde{c}_3 . So if I flip \tilde{c}_3 I should get a legit codeword again:

$$\tilde{c} = (0, 1, 0, 1, 0, 1, 0) \text{ or}$$



And, \tilde{c} is the ONLY solution because flipping any other bit would mess up other circles.

Hooray! That works. But it seems pretty ad hoc. For the rest of this lecture and some of next one, we'll try to introduce some formalism to make this solution seem less ad hoc. At the same time we will flesh out what we care about for ECCs.

these things.

- ① We should be able to handle the **SOMETHING BAD**, whatever that means.
- ② We should be able to recover **WHAT WE WANT TO KNOW** about x .
- ③ We want to **MINIMIZE OVERHEAD**: k/n should be as small as possible.
- ④ We want to do all this **EFFICIENTLY**.

First some definitions:

DEF. The **HAMMING DISTANCE** between $x, y \in \Sigma^n$ is

$$\Delta(x, y) := \sum_{i=1}^n \mathbb{1}\{x_i \neq y_i\}$$

The **RELATIVE HAMMING DISTANCE** between $x, y \in \Sigma^n$ is

$$S(x, y) := \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{x_i \neq y_i\} = \frac{\Delta(x, y)}{n}$$

Notice:
this is a metric.
In particular,
it obeys the
triangle inequality.

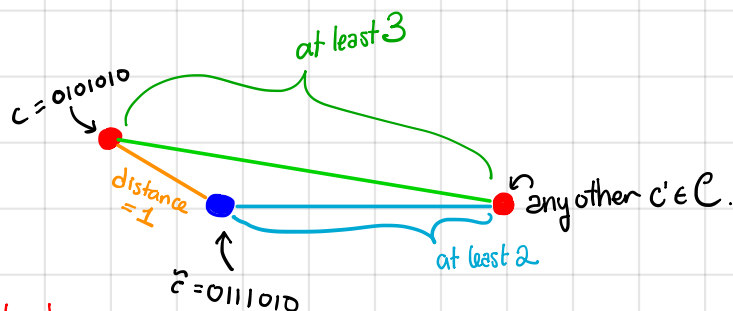
DEF. The **MINIMUM DISTANCE** of a code $C \subseteq \Sigma^n$ is

$$\min_{\substack{c \neq c' \\ c, c' \in C}} \Delta(c, c')$$

Sometimes I'll just
call this "distance."

CLAIM The code in **EXAMPLE 3** has minimum distance 3.

If the CLAIM is true, it explains why that code can correct one error:



DISCLAIMER


I will frequently draw pictures as though Hamming distance is Euclidean distance, and $\{0,1\}^n$ is \mathbb{R}^2 .

Indeed, if $\Delta(c, c') \geq 3 \forall c \neq c' \in \mathcal{C}$, then

$$\Delta(\tilde{c}, c) = 1 \implies \Delta(\tilde{c}, c') \geq 2 \quad \forall c' \in \mathcal{C} \text{ other than } c.$$

by the triangle inequality. Thus, the "correct" codeword $c \in \mathcal{C}$ is uniquely defined by "the one that is closest to \tilde{c} ."

To prove the CLAIM:

- You can probably convince yourself by staring at  (in the same way we convinced ourselves we could always fix one error).
- But we'll see a much less ad hoc way to establish distance after we build up some machinery for **LINEAR CODES** in Lecture 2, so let's put it aside for now.

The POINT of this discussion was that:

MINIMUM DISTANCE is a reasonable proxy for robustness.

aka, $\left\{ \begin{array}{l} \textcircled{1} \text{ We should be able to handle the } \text{SOMETHING BAD}, \text{ whatever that means.} \\ \textcircled{2} \text{ We should be able to recover } \text{WHAT WE WANT TO KNOW} \text{ about } x. \end{array} \right.$

That is,

- In EXAMPLE 2, the code had minimum distance 2 (check this!) and could CORRECT 1 ERASURE and DETECT 1 ERROR.
- In EXAMPLE 3, the code had minimum distance 3, and could CORRECT 1 ERROR.

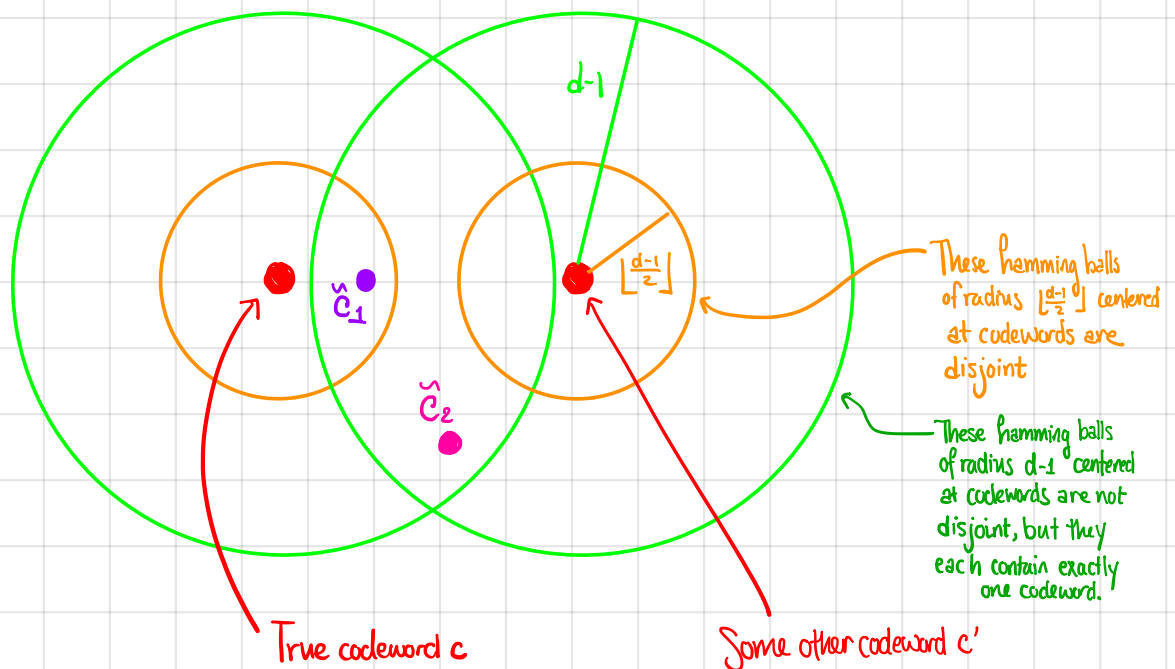
More generally, a code with distance d can:


- correct $\leq d-1$ erasures
- detect $\leq d-1$ errors
- correct $\leq \lfloor \frac{d-1}{2} \rfloor$ errors


For these two, the (inefficient) algorithm is:
"if you see \tilde{c} , return $c \in \mathcal{C}$ that's closest to \tilde{c} ."

For this one, the (inefficient) alg. is:
"If $\tilde{c} \notin \mathcal{C}$, say that something is wrong."

The picture looks like this:



- If c is the "correct" codeword and $\leq \lfloor \frac{d-1}{2} \rfloor$ errors are introduced, we may end up with \tilde{c}_1 . Since all the  balls are disjoint, we can find c from \tilde{c}_1 .

- However, if $\leq d-1$ errors are introduced, we may end up with \tilde{c}_2 . Now it's possible that \tilde{c}_2 came from c or that it came from c' ; we can't tell. However, since each  ball doesn't contain any codeword other than its center, we can tell that something went wrong.

Returning to this, we can now clarify the first two things.

(from earlier)

THINGS WE CARE ABOUT:

- ① We should be able to handle the **SOMETHING BAD**, whatever that means.
- ② We should be able to recover **WHAT WE WANT TO KNOW** about x .
- ③ We want to **MINIMIZE OVERHEAD**: k/n should be as small as possible.
- ④ We want to do all this **EFFICIENTLY**.

If we want:

- ① We should be able to handle $\lfloor \frac{d-1}{2} \rfloor$ **WORST-CASE ERRORS** or $d-1$ **WORST-CASE ERASURES**
- ② We want to recover **ALL OF x** (aka correct the errors or erasures)

Then we should say 1 & 2 We want **MINIMUM DISTANCE d** .

Next we will move on to 3.

ASIDE. A natural question at this point is, "what if I don't want to handle worst-case errors/erasures?" For example, if my code has minimum distance d , and I have two codewords:

$$c = (0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0) \in \{0,1\}^n$$
$$c' = (\underbrace{1 \ 1 \ 1 \ 1}_d \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0) \in \{0,1\}^n$$

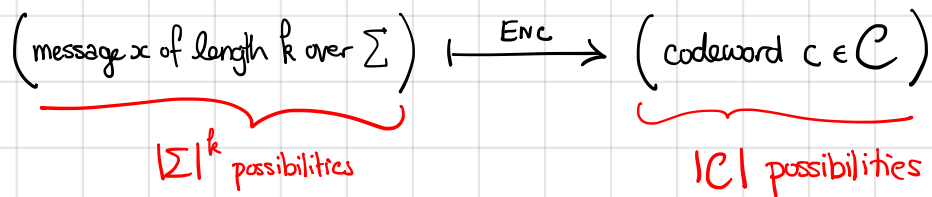
Then if an adversary chooses to flip the first two bits, we'd be confused. But instead say two bits get flipped at random. The probability we get confused is $\frac{\binom{d}{2}}{\binom{n}{2}}$ which might be quite small!

The random-error model (also called the "Shannon model" or "Stochastic model") is natural and important! We will discuss it a little bit in this class. However, most of our focus will be in the worst-case model (also called the "Hamming model" or "adversarial model.")

Moving on to (3), what do we mean by "overhead"?

DEF. The MESSAGE LENGTH (sometimes called DIMENSION) of a code \mathcal{C} over an alphabet Σ is defined to be $k = \log_{|\Sigma|} |\mathcal{C}|$.

This definition makes sense with our operational understanding:



So $|\Sigma|^k = |\mathcal{C}|$ aka $k = \log_{|\Sigma|} |\mathcal{C}|$.

DEF. The RATE of a code $\mathcal{C} \subseteq \Sigma^n$ with block length n over an alphabet Σ is

$$R = \frac{\log_{|\Sigma|} |\mathcal{C}|}{n} = \frac{\text{message length } k}{\text{block length } n}$$

So if R is close to 1, that's GOOD. Not much overhead.

And if R is close to 0, that's BAD. Lots of overhead.

DEF. A code with distance d , message length k , block length n , and alphabet Σ is called a $(n, k, d)_{|\Sigma|}$ code.

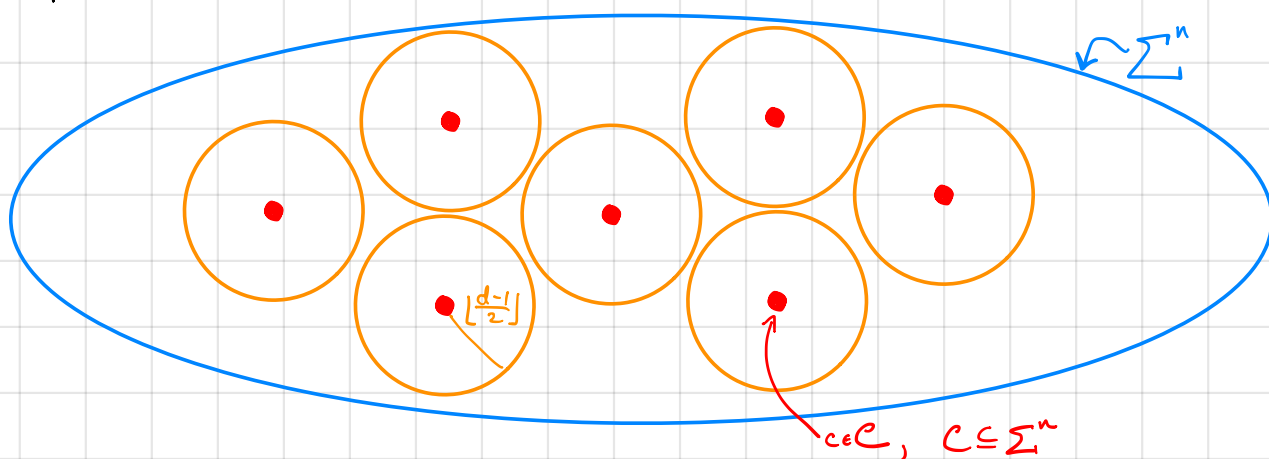
QUESTION. WHAT IS THE BEST TRADE-OFF BETWEEN RATE AND DISTANCE?

This question is still open for binary codes!
But there's lots we do know.

5) RATE vs. DISTANCE: HAMMING BOUND.

What is the best trade-off between rate and distance we can hope for?
The HAMMING BOUND gives one bound on this.

Let's return to the picture we had before, with disjoint Hamming balls of radius $\lfloor \frac{d-1}{2} \rfloor$:



- We have $|\mathcal{C}|$ disjoint Hamming balls of radius $\lfloor \frac{d-1}{2} \rfloor$.
- There can't be too many of them or they wouldn't all fit in Σ^n .

To be a bit more precise:

DEF. The HAMMING BALL in Σ^n of radius e about $x \in \Sigma^n$ is

$$B_{\Sigma^n}(x, e) := \{ y \in \Sigma^n : \Delta(x, y) \leq e \}$$

The VOLUME of $B_{\Sigma^n}(x, e)$ is $\text{Vol}_{|\Sigma|}(e, n) := |B_{\Sigma^n}(x, e)|$

Notice that $|B_{\Sigma^n}(x, e)|$ does not depend on x .

Say that $|\Sigma| = q$. Then

$$\text{Vol}_q(e, n) = 1 + \binom{n}{1}(q-1) + \binom{n}{2}(q-1)^2 + \dots + \binom{n}{e}(q-1)^e$$

\uparrow \uparrow \uparrow \uparrow
 0 all the elements of Σ^n of weight 1 all the elements of Σ^n of weight 2 ... all the elements of Σ^n of weight e .

Notes:

- Sometimes I will drop the " Σ^n " from the $B_{\Sigma^n}(x, e)$ notation
- Sometimes I will write $B_{\Sigma^n}(x, e/n)$ if it's more convenient to talk about relative distance.

So that means that if a code $\mathcal{C} \subseteq \Sigma^n$ has distance d and message length k , where $|\Sigma|=q$,

$$|\mathcal{C}| \cdot \text{Vol}_q \left(\lfloor \frac{d-1}{2} \rfloor, n \right) \leq q^n$$

total volume in the s
total volume in Σ^n

so taking logs of both sides,

$$\log_q(|\mathcal{C}|) + \log_q \left(\text{Vol}_q \left(\lfloor \frac{d-1}{2} \rfloor, n \right) \right) \leq n$$

\swarrow
 $\log_q(|\mathcal{C}|) = k$

$$\Rightarrow \text{Rate} = \frac{k}{n} \leq 1 - \frac{\log_q \left(\text{Vol}_q \left(\lfloor \frac{d-1}{2} \rfloor, n \right) \right)}{n}$$

This is called the **HAMMING BOUND**.

Back to EXAMPLE 3, which was a $(7, 4, 3)_2$ code

\uparrow \uparrow \uparrow \uparrow
 n k d q

- We have $\lfloor \frac{d-1}{2} \rfloor = 1$
- $\text{Vol}_2(1, 7) = 1 + \binom{7}{1} \cdot 1 = 8$
- So

$$\frac{k}{n} \leq 1 - \frac{\log_2(8)}{7} = 1 - \frac{3}{7} = \frac{4}{7}$$

- And in fact $\frac{k}{n} = \frac{4}{7}$, so in this case the Hamming bound is tight!

Notes about this example:

- When the Hamming bound is tight, we say the code is **PERFECT**.
- EXAMPLE 3 (which is perfect) is a special case of something called a **HAMMING CODE**.
- You will explore this more in in-class exercises and on homework.

⑥ RECAP

Now we understand the first 3 of our desiderata:

THINGS WE CARE ABOUT:

- THESE THREE
- ① We should be able to handle the **SOMETHING BAD**, whatever that means.
 - ② We should be able to recover **WHAT WE WANT TO KNOW** about x .
 - ③ We want to **MINIMIZE OVERHEAD**: k/n should be as large as possible.
 - ④ We want to do all this **EFFICIENTLY**.

That is, (for now), our goal is to design codes $C \subseteq \Sigma^n$ so that:

- The **DISTANCE** of C is as large as possible.
- The **RATE** of C is as close to 1 as possible.

Even without the algorithmic considerations, understanding the trade-off between rate and distance turns out to be a fascinating combinatorial question!

In fact, for binary codes ($|\Sigma|=2$), this question is **STILL OPEN!**
(We saw that **EXAMPLE 3** was optimal for $n=7$ and $k=4$, but what about in general?)

Next time, we'll give an overview of abstract algebra, and then give some more definitions that will further de-ad-hocify **EXAMPLE 3**.

That's it for today.

QUESTIONS to PONDER:

- ① How would you generalize the code in EXAMPLE 3 to larger n ?
- ② What is the best bound you can come up with on the rate of a code $C \subseteq \{0,1\}^n$ with distance d ?