## LECTURE 7
## INTRODUCTION TO LINEAR CODES
Information Theory

---

Thomas Debris-Alazard

Inria, École Polytechnique

*Lecture 6:*

*It is possible to communicate Rn bits by sending n bits through a noisy channel Q if and only if*

$$R \leq C(Q) \ \left(\textit{capacity}\right)$$

$\longrightarrow$ The proof relies on the use of **block-codes**:

we encode a symbol into a block-code which adds **redundancy**

An example: spell your name over the phone, send first names!

**M** like Mike, **O** like Oscar, **R** like Romeo, **A** like Alpha, **I** like India and **N** like November

Block-codes reach the capacity of discrete memoryless channels, but. . .

▶ To encode messages to send we need to store a table of exponential size. . .

▶ Encoding is an issue but also decoding, *i.e.,* recovering the sent message from a noisy version

$\Big($in Shannon's proof we need to compute an exponential number of probabilities$\Big)$

Our wish list: defining a sub-class of codes verifying

1. Admitting an efficient encoding algorithm

2. Admitting an efficient decoding algorithm

3. Reaching the capacity

$\longrightarrow$ Linear codes! At least they verify 1...

# BASICS ON LINEAR CODES

A finite field $\mathbb{F}_q$ is a finite set with size $q$ admitting operations $\left(+, -, \times, /\right)$

▶ We necessarily have $q = p^m$ for some prime number $p$ and $m$ integer $> 0$

▶ Algebraic structure: $\mathbb{F}_{p^m} = \mathbb{F}_p[X]/(P(X))$ where $P \in \mathbb{F}_p[X]$ is a polynomial of degree $m$ and irreducible $\left(P = QR, \text{ implies that } P \text{ or } Q \in \mathbb{F}_q\right)$

$$\mathbb{F}_4 = \mathbb{F}_2[X]/(1 + X + X^2)$$
$$X(1 + X) = X + X^2 = -1 = 1$$

**Be careful:**

$$\mathbb{F}_q = \mathbb{Z}/q\mathbb{Z} \iff q \text{ is prime}$$

**An important example: the binary field $\mathbb{F}_2$**

$$\mathbb{F}_2 = \{0, 1\} \text{ where}$$
$$0 + 1 = 1 + 0 = 1, 0 + 0 = 1 + 1 = 0, 1 \times 0 = 0 \times 1 = 0 \times 0 = 0 \text{ and } 1 \times 1 = 1$$

$$\mathbb{F}_q^n = \underbrace{\mathbb{F}_q \times \cdots \times \mathbb{F}_q}_{n \text{ times}} \text{ is a } \mathbb{F}_q\text{-vector space}$$

$$(x_1, \ldots, x_n) + (y_1, \ldots, y_n) = (x_1 + y_1, \ldots, x_n + y_n)$$

$$\forall \lambda \in \mathbb{F}_q, \ \ \lambda \cdot (x_1, \ldots, x_n) = (\lambda x_1, \ldots, \lambda x_n)$$

**Linear codes:**

A linear code $\mathcal{C}$ is a subspace of $\mathbb{F}_q^n$

When $\mathcal{C}$ has dimension $k$, we say that it is an $[n, k]_q$-code: $n$ length, $k$ dimension

*Linear codes are block-codes when the alphabet is a finite field $+$ a linear structure $\left(\text{subspace}\right)$*

**First example: repetition code of length 3**

$$\left\{(0, 0, 0), (1, 1, 1)\right\} \text{ is a } [3, 1]_2\text{-code}$$

**Rate of linear codes:**

An $[n, k]_q$-code has cardinal $q^k$ $\left(\text{why?}\right)$ and its rate $R$ is equal to

$$R = \frac{\log_q q^k}{n} = \frac{k}{n}$$

**Non trivial linear codes:**

1. $\left\{ (f(x_1), \ldots, f(x_n)) : f \in \mathbb{F}_q[X], \ \deg(f) < k \right\} \subseteq \mathbb{F}_q^n$

2. Given two linear codes $U, V \subseteq \mathbb{F}_q^{n/2}$, $\left\{ (\mathbf{u}, \mathbf{u} + \mathbf{v}) : \mathbf{u} \in U \text{ and } \mathbf{v} \in V \right\} \subseteq \mathbb{F}_q^n$

**Exercise Session:**

What are the dimensions of the above linear codes?

*How to represent an $[n, k]_q$-code? It has size $q^k$, is a table of this size necessary?*

*How to represent an $[n, k]_q$-code? It has size $q^k$, is a table of this size necessary?*

**No!**

**Basis / Primal representation:**

An $[n, k]_q$-code $\mathcal{C}$ admits a basis $\mathbf{b}_1, \ldots, \mathbf{b}_k \in \mathbb{F}_q^n$

$$\mathcal{C} = \left\{ \mathbf{m}\mathbf{G} : \mathbf{m} \in \mathbb{F}_q^k \right\} \text{ where the rows of } \mathbf{G} \in \mathbb{F}_q^{k \times n} \text{ are the } \mathbf{b}_i\text{'s}$$

The matrix $\mathbf{G}$ is called a **generator matrix** of $\mathcal{C}$

**Redundancy versus rate:**

Given a binary code $\mathcal{C}$ of dimension $k$, we can **easily** encode $k$ bits $(m_1, \ldots, m_k)$ as $\mathbf{m}\mathbf{G} \in \mathcal{C}$ where

$\mathbf{G}$ generator matrix $\Big($**to encode does not necessitate to store an exponential size table**$\Big)$
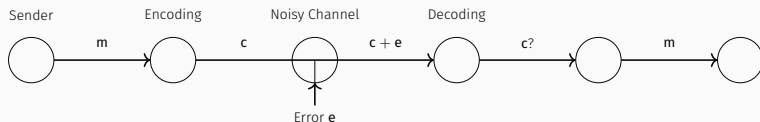
We have mapped $k$ bits to $n$ bits! The $\Big($normalized$\Big)$ redundancy $(n - k)/n = 1 - k/n = 1 - R$

$R \approx 0$: a lot of redundancy  ;  $R \approx 1$: few redundancy

*Particular case of $\mathbb{F}_2$, but can be generalized to $\mathbb{F}_q$ by encoding elements with $\log_2(q)$ bits*

How to transmit $k$ bits over a noisy channel?

1. **Linear code:** fix $\mathcal{C}$ subspace $\subseteq \mathbb{F}_2^n$ of dimension $k < n$

2. **Encoding:** map $(m_1, \ldots, m_k) \longrightarrow c = (c_1, \ldots, c_n) \in \mathcal{C}$ task adding $n - k$ bits redundancy

   $\longrightarrow$ as $\mathcal{C}$ is linear the encoding is easy $\Big($ only linear algebra $\Big)$, *i.e.,* $c = mG$

3. Send $c$ across the noisy channel, errors happen and some bits of $c$ are modified



Sender — Encoding — Noisy Channel — Decoding

$m$ — $c$ — $c + e$ — $c?$ — $m$

Error $e$

**Decoding:**

$\longrightarrow$ from $c + e$: recover $e$ and then $c$. Now as $G$ has rank $k$, we easily recover $m$

by Gaussian elimination $\Big($ we use the linearity $\Big)$

# DUAL REPRESENTATION OF CODES

*Linear codes as subspaces can also be written as the kernel of a matrix*

**Dual code:**

Given an $[n, k]_q$-code $\mathcal{C}$, its dual $\mathcal{C}^\perp$ is an $[n, n-k]_q$-code defined as

$$\mathcal{C}^\perp = \left\{ \mathbf{c}^\perp \in \mathbb{F}_q^n : \ \forall \mathbf{c} \in \mathcal{C}, \ \langle \mathbf{c}^\perp, \mathbf{c} \rangle \overset{\text{def}}{=} \sum_{i=1}^n \underbrace{c_i^\perp c_i}_{\in \mathbb{F}_q} = 0 \right\}$$

**Parity-check/Dual representation:**

$\mathcal{C}^\perp$ is an $[n, n-k]_q$-code. Furthermore, for any generator matrix $\mathbf{H} \in \mathbb{F}_q^{(n-k) \times n}$ $\Big($rows of $\mathbf{H}$

form a basis of $\mathcal{C}^\perp\Big)$ we have,

$$\mathcal{C} = \left\{ \mathbf{c} \in \mathbb{F}_q^n : \ \mathbf{H}\mathbf{c}^\top = \mathbf{0} \right\}$$

Furthermore, any matrix $\mathbf{H} \in \mathbb{F}_q^{(n-k) \times n}$ with rank $n-k$, such that $\mathcal{C}$ is its right kernel, forms

$\Big($considering its rows$\Big)$ a basis of $\mathcal{C}^\perp$.

Such matrix $\mathbf{H}$ is called a parity-check matrix of $\mathcal{C}$

**Proof:**

It is clear that $\mathcal{C}^{\perp}$ is a subspace of $\mathbb{F}_q^n$. Let us show that $\mathcal{C}$ has dimension $n - k$. First, $\mathcal{C}$ can be written as the right kernel of a matrix $\mathbf{H} \in \mathbb{F}_q^{(n-k)\times n}$ with rank $n - k$,

$$\mathcal{C} = \left\{ \mathbf{c} \in \mathbb{F}_q^n : \ \mathbf{H}\mathbf{c}^{\top} = \mathbf{0} \right\}$$

Therefore, all rows of $\mathbf{H}$ are elements in $\mathcal{C}^{\perp}$ showing that $\dim \mathcal{C}^{\perp} \geq n - k$. On the other hand, if $\mathbf{B} \in \mathbb{F}_q^{m\times n}$ is a basis $\left(\text{considering its rows}\right)$ of $\mathcal{C}^{\perp}$. Then by linearity $\mathcal{C}$ is included in the $\left(\text{right}\right)$ kernel of $\mathbf{B}$. We deduce that $k = \dim \mathcal{C} \leq n - \dim \mathcal{C}^{\perp}$ concluding the whole proof

$\mathsf{G} \in \mathbb{F}_q^{k \times n}$ generator matrix of $\mathcal{C}$, *i.e.,* $\mathcal{C} = \left\{ \mathsf{mG} : \ \mathsf{m} \in \mathbb{F}_q^k \right\}$

$\longrightarrow$ **SG** is still a generator matrix when $\mathsf{S} \in \mathbb{F}_q^{k \times k}$ is invertible

$\mathsf{H} \in \mathbb{F}_q^{(n-k) \times n}$ parity-check matrix of $\mathcal{C}$, *i.e.,* $\mathcal{C} = \left\{ \mathsf{c} \in \mathbb{F}_q^n : \ \mathsf{Hc}^\top = \mathsf{0} \right\}$

$\longrightarrow$ **SH** is still a parity-check matrix when $\mathsf{S} \in \mathbb{F}_q^{(n-k) \times (n-k)}$ is invertible

*Left multiplication by an invertible matrix computes a change of basis!*

$\mathsf{G} \in \mathbb{F}_q^{k \times n}$ generator matrix of $\mathcal{C}$ $\xleftrightarrow{\text{easy to compute?}}$ $\mathsf{H} \in \mathbb{F}_q^{(n-k) \times n}$ parity-check matrix of $\mathcal{C}$

▶ Given $\mathsf{G} \in \mathbb{F}_q^{k \times n}$, it has rank $k$. We can perform a Gaussian elimination, *i.e.*, compute $\mathsf{S} \in \mathbb{F}_q^{k \times k}$ invertible such that $\Big($up to a permutation of columns$\Big)$,

$$\mathsf{SG} = (\mathsf{I}_k \mid \mathsf{A}) \text{ where } \mathsf{A} \in \mathbb{F}_q^{k \times (n-k)}$$

$\longrightarrow$ Then $\mathsf{H} = (-\mathsf{A}^\top \mid \mathsf{I}_{n-k})$ parity-check matrix of $\mathcal{C}$

Proof:

Indeed, $\mathsf{m}\,(\mathsf{SG})\,\mathsf{H}^\top = \mathsf{m}(\mathsf{I}_k \mid \mathsf{A}) \left( \dfrac{-\mathsf{A}}{\mathsf{I}_{n-k}} \right) = \mathsf{m}(0) = (0)$. Therefore, $\mathcal{C}$ included in the right

kernel of $\mathsf{H}$. But $\mathsf{H}$ has rank $n - k$, showing the $\mathsf{H}$ is a parity-check matrix of $\mathcal{C}$

$\mathsf{G} \in \mathbb{F}_q^{k \times n}$ generator matrix of $\mathcal{C}$ $\xleftrightarrow{\text{easy to compute?}}$ $\mathsf{H} \in \mathbb{F}_q^{(n-k) \times n}$ parity-check matrix of $\mathcal{C}$

▶ Given $\mathsf{H} \in \mathbb{F}_q^{(n-k) \times n}$, it has rank $n - k$. We can perform a Gaussian elimination, *i.e.,* compute $\mathsf{S} \in \mathbb{F}_q^{(n-k) \times (n-k)}$ invertible such that $\Big($ up to a permutation of columns $\Big)$,

$$\mathsf{SH} = (\mathsf{I}_{n-k} \mid \mathsf{B}) \text{ where } \mathsf{B} \in \mathbb{F}_q^{(n-k) \times k}$$

$$\longrightarrow \text{Then } \mathsf{G} = (-\mathsf{B}^\top \mid \mathsf{I}_k) \text{ generator matrix of } \mathcal{C}$$

### Exercise:

Given $\mathbf{x} \in \mathbb{F}_q^n$ and a linear code $\mathcal{C} \subseteq \mathbb{F}_q^n$, is it easy to decide if $\mathbf{x} \in \mathcal{C}$?

$$\mathcal{C}^\perp = \left\{ \mathbf{c}^\perp \in \mathcal{C}^\perp : \ \forall \mathbf{c} \in \mathcal{C}, \ \langle \mathbf{c}, \mathbf{c}^\perp \rangle = \sum_{i=1}^{n} c_i^\perp c_i = 0 \in \mathbb{F}_q \right\}$$

If $\mathcal{C} \subseteq \mathbb{F}_q^n$ has dimension $k$, then $\mathcal{C}^\perp$ has dimension $n - k$ where $n = \dim \mathbb{F}_q^n$

$\longrightarrow$ It seems that $\mathcal{C}^\perp$ is the orthogonal of $\mathcal{C}$ and $\langle \cdot, \cdot \rangle$ is a scalar product, but no!

The dual is not an orthogonal!

$$\mathcal{C} + \mathcal{C}^\perp \neq \mathbb{F}_q^n$$

It happens that $\mathcal{C} \cap \mathcal{C}^\perp \neq \{0\}$ and this intersection is called the hull

**Characters and Fourier transforms for prime $q$:**

▶ Characters are $\chi_{\mathbf{x}}(\mathbf{y}) = e^{2i\pi \langle \mathbf{x}, \mathbf{y} \rangle / q}$ where $\mathbf{x}, \mathbf{y} \in \mathbb{F}_q^n$. They are morphisms from $(\mathbb{F}_q^n, +)$ to the units of $(\mathbb{C}, \times)$.

▶ The Fourier transform of $f : \mathbb{F}_q^n \to \mathbb{C}$, is

$$\widehat{f}(\mathbf{x}) = \frac{1}{\sqrt{q}} \sum_{\mathbf{y} \in \mathbb{F}_q^n} f(\mathbf{y}) \chi_{\mathbf{x}}(\mathbf{y})$$

*The dual code is defined via group theory involving dual groups*

The dual code $\mathcal{C}^\perp$ is the set of points for which characters are trivial when restricted to $\mathcal{C}$, *i.e.,*

$$\mathcal{C}^\perp = \left\{ \mathbf{c}^\perp \in \mathbb{F}_q^n : \ \forall \mathbf{c} \in \mathcal{C}, \ \chi_{\mathbf{c}^\perp}(\mathbf{c}) = 1 \right\}$$

$$\left( \text{when } q \text{ prime}, \chi_{\mathbf{x}}(\mathbf{y}) = 1 \Longleftrightarrow \langle \mathbf{x}, \mathbf{y} \rangle = 0 \in \mathbb{F}_q \right)$$

Given two finite subspaces: $F \subseteq E$

Equivalence relation: $x \sim y \iff x - y \in F$

$E/F = \{\bar{x} \ : \ x \in E\}$ where $\bar{x} \stackrel{\text{def}}{=} \{y \in E \ : \ x \sim y\} = x + F$

$\longrightarrow$ It defines a linear space!

$k = \dim E/F = \dim E - \dim F$

Rough analogy:

| $E/F$ | $\mathbb{Z}/4\mathbb{Z}$ |
|---|---|
| $\{\overline{x_1}, \ldots, \overline{x_N}\}$ | $\{\bar{0}, \bar{1}, \bar{2}, \bar{3}\}$ |
| $\overline{x_i} = x_i + F$ | $\bar{\ell} = \ell + 4\mathbb{Z}$ |
| $\bar{x} = \bar{y} \iff x - y \in F$ | $\bar{\ell} = \bar{m} \iff \ell - m \in 4\mathbb{Z}$ |
| $E = \bigsqcup\limits_{1 \leq i \leq N} \overline{x_i}$ | $\mathbb{Z} = \bigsqcup\limits_{\ell \in \{0,1,2,3\}} \bar{\ell}$ |

Decoding: given $\mathbf{c} + \mathbf{e}$, recover $\mathbf{e}$

$\longrightarrow$ Make modulo $\mathcal{C}$ to extract the information about $\mathbf{e}$

**Coset space:** $\mathbb{F}_2^n / \mathcal{C}$

$$\sharp \, \mathbb{F}_q^n / \mathcal{C} = q^{n-k} \quad \text{and} \quad \mathbb{F}_q^n / \mathcal{C} = \left\{ \overline{\mathbf{x}}_i \; : 1 \leq i \leq q^{n-k} \right\} = \left\{ \mathbf{x}_i + \mathcal{C} \; : \; 1 \leq i \leq q^{n-k} \right\}$$

where the $\mathbf{x}_i$'s are the representatives of $\mathbb{F}_q^n / \mathcal{C}$. The $x_i + \mathcal{C}$'s are disjoint!

A natural set of representatives via a parity-check $\mathbf{H}$: syndromes

**Proposition:**

We have:

1. $\mathbf{x}_i + \mathcal{C} \in \mathbb{F}_q^n / \mathcal{C} \longmapsto \mathbf{H}\mathbf{x}_i^{\mathsf{T}} \in \mathbb{F}_q^{n-k}$ $\left( \text{called a syndrome} \right)$ is an isomorphism

2. $\mathbb{F}_q^n = \bigsqcup_{s \in \mathbb{F}_q^{n-k}} \left\{ \mathbf{z} \in \mathbb{F}_q^n \; : \; \mathbf{H}\mathbf{z}^{\mathsf{T}} = \mathbf{s}^{\mathsf{T}} \right\}$

$\mathbf{c} + \mathbf{e} \bmod \mathcal{C} = \mathbf{H}(\mathbf{c}+\mathbf{e})^{\mathsf{T}} = \underbrace{\mathbf{H}\mathbf{c}^{\mathsf{T}}}_{=0} + \mathbf{H}\mathbf{e}^{\mathsf{T}} = \mathbf{H}\mathbf{e}^{\mathsf{T}}$ which gives information to recover $\mathbf{e}$ $\left( \text{decoding} \right)$

$\longrightarrow \mathbf{c} + \mathbf{e} \bmod \mathcal{C}$ is only function of $\mathbf{e}$!

**Proof:**

1. Let us first show that $\mathbf{x}_i + \mathcal{C} \in \mathbb{F}_q^n/\mathcal{C} \longmapsto \mathbf{Hx}_i^{\mathsf{T}} \in \mathbb{F}_q^{n-k}$ is a well-defined mapping. If we choose another class representative $\mathbf{y}_i + \mathcal{C} = \mathbf{x}_i + \mathcal{C}$. Then by definition

$$\mathbf{y}_i - \mathbf{x}_i \in \mathcal{C} \iff \mathbf{H}\left(\mathbf{y}_i - \mathbf{x}_i\right)^{\top} = \mathbf{0} \iff \mathbf{Hy}_i^{\top} = \mathbf{Hx}_i^{\top}$$

It shows that we have a well-defined mapping. But the equivalence also shows that it is a one-to-one mapping

The above application is surjective as $\mathbf{H}$ has rank $n - k$, therefore for any $\mathbf{s} \in \mathbb{F}_q^{n-k}$ it exists $\mathbf{x} \in \mathbb{F}_q^n$ such that $\mathbf{Hx}^{\top} = \mathbf{s}^{\top}$ and $\mathbf{x}$ defines one representative. Furthermore the mapping is clearly linear, concluding the proof of 1

2. This is a consequence of the equivalence relation but let's give a direct proof. We have shown above that $\forall \mathbf{z} \in \mathbb{F}_q^n$, it exists $\mathbf{s} \in \mathbb{F}_q^n$ such that $\mathbf{Hz}^{\top} = \mathbf{s}^{\top}$ $\left(\mathbf{H} \text{ has rank } n - k\right)$.

To conclude notice that $\left\{\mathbf{z} \in \mathbb{F}_q^n \ : \ \mathbf{Hz}^{\mathsf{T}} = \mathbf{s}^{\mathsf{T}}\right\}$ are clearly disjoint for $\mathbf{s} \in \mathbb{F}_q^{n-k}$

$\mathcal{C}$ be an $[n, k]_q$-code with generator and parity-check matrices $\mathsf{G}$ and $\mathsf{H}$

▶ Given a noisy codeword, $\mathbf{y} = \underbrace{\mathbf{c}}_{\in \mathcal{C}} + \mathbf{e}$, its syndrome is

$$\mathsf{H}\mathbf{y}^\top = \mathsf{H}\mathbf{c}^\top + \mathsf{H}\mathbf{e}^\top = \mathsf{H}\mathbf{e}^\top \text{ where we use } \mathcal{C} = \left\{ \mathbf{c} \in \mathbb{F}_q^n : \ \mathsf{H}\mathbf{c}^\top = \mathbf{0} \right\}$$

▶ Given a syndrome, $\mathbf{s}^\top = \mathsf{H}\mathbf{e}^\top$, we can easily compute its associated noisy codeword, by a Gaussian elimination we compute $\mathbf{y}$ such that $\mathsf{H}\mathbf{y}^\top = \mathbf{s}^\top$ $\left( \text{as } \mathsf{rank}(\mathsf{H}) = n - k \right)$

$$\mathsf{H}\mathbf{y}^\top = \mathbf{s}^\top \iff \mathsf{H}(\mathbf{y} - \mathbf{e})^\top = \mathbf{0} \iff \mathbf{y} - \mathbf{e} \in \mathcal{C} \iff \mathbf{y} = \underbrace{\mathbf{c}}_{\in \mathcal{C}} + \mathbf{e}$$

# HAMMING DISTANCE

*Remember, we introduced codes to communicate over a noisy channel*

$\longrightarrow$ We will restrict our attention to the following channels:
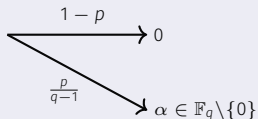
$q$-**ary symmetric channels:**

Memoryless channel $\left( \mathbb{F}_q, \mathbb{F}_q, p(y \mid x) \right)$ where,

$$\forall x, y \in \mathbb{F}_q, \quad p(y \mid x) = \begin{cases} 1 - p & \text{if } x = y \\ \frac{p}{q-1} & \text{otherwise} \end{cases}$$

$p$: probability of error ; $\frac{p}{q-1}$ transition probability

When sending $\mathbf{x} \in \mathbb{F}_q^n$ through the channel

$\mathbf{y} = \mathbf{c} + \mathbf{e}$ where the $e_i$ are i.i.d and $p(e_i = x) = \begin{cases} 1 - p & \text{if } x = 0 \\ \frac{p}{q-1} & \text{otherwise} \end{cases}$

$$\xrightarrow{1-p} 0$$

$$\xrightarrow[\frac{p}{q-1}]{} \alpha \in \mathbb{F}_q \backslash \{0\}$$

*Remember, after sending a codeword across a noisy channel we want to recover the sent codeword, i.e., decoding*

$\longrightarrow$ The optimal decoder is the following one:

**Maximum likelihood decoder:**

Given a $q$-ary symmetric channel $\left(\mathbb{F}_q, \mathbb{F}_q, p(y \mid x)\right)$ and a block-code $\mathcal{C} \subseteq \mathbb{F}_q^n$. We call the maximum likelihood decoder the map

$$\varphi : \mathbb{F}_q^n \longrightarrow \mathcal{C}$$

such that given $\mathbf{y} \in \mathbb{F}_q^n$, it outputs the codeword $\mathbf{c} \in \mathcal{C}$ maximizing the transition probabilities

$$\varphi(\mathbf{y}) \stackrel{\text{def}}{=} \arg\max_{\mathbf{c} \in \mathcal{C}} p(\mathbf{c} \mid \mathbf{y}) = \arg\max_{\mathbf{c} \in \mathcal{C}} \prod_{i=1}^n p(c_i \mid y_i)$$

**Proposition:**

In a $q$-ary symmetric channel with probability of transition $p/(q-1) < 1/q$, if codewords $\mathbf{c} \in \mathcal{C}$ are chosen uniformly at random among $\mathcal{C}$, then

$$\forall \mathbf{y} \in \mathbb{F}_q^n, \ \varphi(\mathbf{y}) = \mathbf{c} \in \mathcal{C} \ \text{ such that } \mathbf{c} = \underset{\mathbf{d} \in \mathcal{C}}{\arg \min} \, d_H(\mathbf{y}, \mathbf{d})$$

where $d_H(\cdot, \cdot)$ is the Hamming distance,

$$\forall \mathbf{x}, \mathbf{y} \in \mathbb{F}_q^n, \ d_H(\mathbf{x}, \mathbf{y}) \overset{\text{def}}{=} \sharp \{i \in [1, n], \ x_i \neq y_i\}$$

Given $\mathbf{y} \in \mathbb{F}_q^n$, the maximum likelihood decoder for $q$-ary symmetric channels outputs the closest codewords $\mathbf{c}$ for the Hamming distance

It justifies the use of the Hamming distance for decoding in the $q$-ary symmetric channel

**Proof:**

1. First, by using that the **c**'s are uniform among $\mathcal{C}$

$$p(\mathbf{c} \mid \mathbf{y}) = p(\mathbf{y} \mid \mathbf{c}) \frac{p(\mathbf{c})}{p(\mathbf{y})} = p(\mathbf{y} \mid \mathbf{c}) \frac{1}{\sharp \mathcal{C} \, p(\mathbf{y})}$$

   We deduce that maximizing $p(\mathbf{c} \mid \mathbf{y})$ $\left(\text{over } \mathbf{c}\right)$ boils down to maximize $p(\mathbf{y} \mid \mathbf{c})$

2. Second, by definition of the $q$-ary symmetric channel,

$$p(\mathbf{y} \mid \mathbf{c}) = (1 - p)^{\sharp\{i \in [1,n]:\, y_i = c_i\}} \left(\frac{p}{q-1}\right)^{\sharp\{i \in [1,n]:\, y_i \neq c_i\}}$$

$$= (1 - p)^{n - d_H(\mathbf{y}, \mathbf{c})} \left(\frac{p}{q-1}\right)^{d_H(\mathbf{y}, \mathbf{c})}$$

   As $p/(q-1) < 1/q$,

$$\alpha \mapsto (1 - p)^{n - \alpha} \left(\frac{p}{q-1}\right)^{\alpha}$$

   is a decreasing function showing that $p(\mathbf{y} \mid \mathbf{c})$ is maximal when $d_H(\mathbf{y}, \mathbf{c})$ is minimal

**Hamming weight:**

$$\forall \mathbf{x} \in \mathbb{F}_q^n, \ |\mathbf{x}| \stackrel{\text{def}}{=} \sharp \{i \in [1, n], \ x_i \neq 0\}$$

$$\longrightarrow \ d_H(\mathbf{x}, \mathbf{y}) = |\mathbf{x} - \mathbf{y}|$$

**Some remarks:**

- $|\cdot|$ is not a norm but $d_H(\cdot, \cdot)$ is a distance

- The Hamming weight does not discriminate non-zero symbols, for instance in $\mathbb{F}_5 = \mathbb{Z}/5\mathbb{Z}$,
$$\left| (1, 2, 0, 1, 0, 0, 2) \right| = \left| (3, 3, 4, 0, 0, 0, 1) \right| = 4$$

An important parameter for a code: its minimum distance

$\longrightarrow$ It measures the quality of a code in terms of "error detection"

**Minimum distance:**

Given $\mathcal{C} \subseteq \mathbb{F}_q^n$, its minimum distance is defined as

$$d_{\min}(\mathcal{C}) \stackrel{\text{def}}{=} \min \left\{ |c_1 - c_2| : \ c_1, c_2 \in \mathcal{C} \text{ and } c_1 \neq c_2 \right\}$$

**Remark:**

For a linear code $\mathcal{C}$,

$$d_{\min}(\mathcal{C}) = \min \left\{ |c| : \ c \in \mathcal{C} \setminus \{0\} \right\}$$

Suppose that someone sends us a codeword $c \in \mathcal{C}$ across a noisy channel

Our goal is to guess if an error occurred

*How can we proceed? What is the maximal amount of errors for which we can take the right decision with certainty?*

Suppose that someone sends us a codeword $\mathbf{c} \in \mathcal{C}$ across a noisy channel

Our goal is to guess if an error occurred

*How can we proceed? What is the maximal amount of errors for which we can take the right decision with certainty?*

**Error detection strategy:**

Given a received $\mathbf{y}$ we compute $\mathbf{H}\mathbf{y}^\top$ for $\mathbf{H}$ being a parity-check matrix of the code. If we obtain $\mathbf{0}$ then we say that no error occurred

This strategy gives the right answer with certainty if the Hamming weight of the error is $< d_{\min}(\mathcal{C})$!

**Proof:**

If an error occurred then we receive $\mathbf{c} + \mathbf{e}$. Therefore $\mathbf{H}\left(\mathbf{c} + \mathbf{e}\right)^\top = \mathbf{H}\mathbf{c}^\top + \mathbf{H}\mathbf{e}^\top = \mathbf{H}\mathbf{e}^\top$. Then if $|\mathbf{e}| < d_{\min}(\mathcal{C})$ we necessarily have $\mathbf{e} \notin \mathcal{C}$ and $\mathbf{H}\mathbf{e}^\top \neq \mathbf{0}$. However, if $|\mathbf{e}| \geq d_{\min}(\mathcal{C})$ it is possible that $\mathbf{e} \in \mathcal{C}$ and $\mathbf{H}\mathbf{e}^\top = \mathbf{0}$

$$\mathbf{x} \in \mathbb{F}_q^n, \quad \mathcal{B}(\mathbf{x}, r) \stackrel{\text{def}}{=} \left\{ \mathbf{y} \in \mathbb{F}_q^n : |\mathbf{y} - \mathbf{x}| \leq r \right\}$$

**Proposition:**

Given a code $\mathcal{C} \subseteq \mathbb{F}_q^n$,

$$\forall \mathbf{c}_1, \mathbf{c}_2 \in \mathcal{C}, \quad \mathbf{c}_1 \neq \mathbf{c}_2: \quad \mathcal{B}\left(\mathbf{c}_1, \left\lfloor \frac{d_{\min}(\mathcal{C})-1}{2} \right\rfloor\right) \bigcap \mathcal{B}\left(\mathbf{c}_2, \left\lfloor \frac{d_{\min}(\mathcal{C})-1}{2} \right\rfloor\right) = \emptyset$$

**Proof:**

By contradiction, suppose there exists $\mathbf{y} \in \mathcal{B}\left(\mathbf{c}_1, \left\lfloor \frac{d_{\min}(\mathcal{C})-1}{2} \right\rfloor\right) \bigcap \mathcal{B}\left(\mathbf{c}_2, \left\lfloor \frac{d_{\min}(\mathcal{C})-1}{2} \right\rfloor\right)$,

$$\begin{aligned}
|\mathbf{c}_1 - \mathbf{c}_2| = |(\mathbf{c}_1 - \mathbf{y}) - (\mathbf{c}_2 - \mathbf{y})| \\
\leq |\mathbf{c}_1 - \mathbf{y}| + |\mathbf{c}_2 - \mathbf{y}| \quad \text{(triangular inequality)} \\
\leq \left\lfloor \frac{d_{\min}(\mathcal{C}) - 1}{2} \right\rfloor + \left\lfloor \frac{d_{\min}(\mathcal{C}) - 1}{2} \right\rfloor \\
< d_{\min}(\mathcal{C})
\end{aligned}$$

which is a contradiction as $\mathbf{c}_1 \neq \mathbf{c}_2$ and they belong to $\mathcal{C}$ with minimum distance $d_{\min}(\mathcal{C})$

When transmitting $\mathbf{c} \in \mathcal{C}$, if the Hamming weight of the error is $< d_{\min}(\mathcal{C})/2$, then the maximum likelihood decoder necessarily outputs $\mathbf{c}$ $\left(\text{closest codeword for the Hamming distance}\right)$

When transmitting $\mathbf{c} \in \mathcal{C}$, if the Hamming weight of the error is $< d_{\min}(\mathcal{C})/2$, then the maximum

likelihood decoder necessarily outputs $\mathbf{c}$

The above statement says that with $< d_{\min}(\mathcal{C})/2$ error the maximum likelihood decoder succeeds

with certainty!

$\longrightarrow$ There are codes for which the maximum likelihood decoder works with probability $1 - e^{-Cn}$

as soon as there are $\leq d_{\min}(\mathcal{C})$ errors, we gain a factor two!

$\Big($in particular random codes as we will see in Lecture 8$\Big)$

*The minimum distance quantifies how "good" is a code in term of error decoding/detection*

▶ Balls centered at codewords with radius $\left\lfloor \frac{d_{\min}(\mathcal{C})-1}{2} \right\rfloor$ are disjoint

$\longrightarrow$ We can correct $\left\lfloor \frac{d_{\min}(\mathcal{C})-1}{2} \right\rfloor$ errors with certainty!

▶ There are no codewords in any ball centered at codewords with radius $d_{\min}(\mathcal{C}) - 1$

$\longrightarrow$ We can detect any $< d_{\min}(\mathcal{C})$ errors

Given an $[n, k]_q$-code $\mathcal{C}$, how large can be its minimum distance $d_{\min}(\mathcal{C})$?

# BOUNDS ON MINIMUM DISTANCE

The $[7, 4]_2$ Hamming code $\mathcal{C}_H$ admits as parity check matrix

$$\mathsf{H} \stackrel{\text{def}}{=} \begin{pmatrix} 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ 0 & 1 & 1 & 0 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 & 1 & 0 & 1 \end{pmatrix}$$

It has minimum distance $d_{\min}(\mathcal{C}_H) = 2$, indeed use the following proposition

**Proposition:**

Given a linear code $\mathcal{C}$ with parity-check matrix $\mathsf{H}$,

$$\mathcal{C} \text{ has minimum distance} \geq d \iff \text{every } d - 1 \text{ columns of } \mathsf{H} \text{ form a free family}$$

Given $\mathbf{c} \in \mathcal{C}_H$ there are $8 = 2^3$ noisy codewords $\mathbf{c} + \mathbf{e}$ where $|\mathbf{e}| \leq 1 = \left\lfloor \frac{d_{\min}(\mathcal{C}_H) - 1}{2} \right\rfloor$

$$\longrightarrow \text{The balls } \mathcal{B}\left(\mathbf{c}, \left\lfloor \frac{d_{\min}(\mathcal{C}) - 1}{2} \right\rfloor\right)\text{'s for } \mathbf{c} \in \mathcal{C}_H \text{ form a partition of } \mathbb{F}_2^7!$$

**Perfect codes:**

A linear code $\mathcal{C} \subseteq \mathbb{F}_q^n$ is said to be perfect if the balls $\mathcal{B}\left(\mathbf{c}, \left\lfloor \frac{d_{\min}(\mathcal{C}) - 1}{2} \right\rfloor\right)$ form a partition of $\mathbb{F}_q^n$

**Exercise:**

Given a parity-check of some matrix $H$, is it easy to check that every $d - 1$ columns of $H$ form a free family? More generally, does it seem easy to compute the minimum distance of a given code?

An Hamming code is the $[2^m - 1, 2^m - m - 1]_2$-code admitting as parity-check matrix

$H \in \mathbb{F}_2^{(2^m-1) \times m}$ whose columns are all the vectors $\mathbb{F}_2^m \setminus \{\mathbf{0}\}$. It has minimum distance 3 and it is a

perfect code

$$2^{2^m - m - 1} \left( \binom{2^m - 1}{1} + 1 \right) = 2^{2^m - 1}$$

**Theorem:**

Parameters $[n, k, d_{min}(\mathcal{C})]_q$ of perfect codes are known: $[2\ell + 1, 1, 2\ell + 1]_2$ $\left(\text{repetition codes with}\right.$

odd length$\left.\right)$, $[2^m - 1, 2^m - m - 1, 3]_2$ $\left(\text{Hamming codes}\right)$, $[23, 12, 7]_2$ $\left(\text{binary Golay code } G_{23}\right)$

and $[11, 6, 5]_3$ $\left(\text{ternary Golay code } G_{11}\right)$

**Singleton bound:**

For all $[n, k]_q$-code $\mathcal{C}$,

$$d_{\min}(\mathcal{C}) \le n - k + 1$$

**Proof:**

Given a parity-check matrix $\mathbf{H} \in \mathbb{F}_q^{(n-k) \times n}$, it has rank $n - k$. We cannot hope having more than $(n - k)$ columns forming a free family. Therefore,

$$d - 1 \le n - k$$

*Do we know codes that reach this bound?* Yes!

$\Big($codes reaching the Singleton bound are said MDS, *i.e.,* Maximum Distance Separable$\Big)$

**Reed-Solomon code:**

Given $x_1, \ldots, x_n \in \mathbb{F}_q$ where the $x_i$'s are different, *i.e.*, $x_i \neq x_j$,

$$\left\{ (f(x_1), \ldots, f(x_n)) : f \in \mathbb{F}_q[X], \ \deg(f) < k \right\}$$

is a $[n, k]_q$-code with minimum distance $n - k + 1$

Reed-Solomon codes have optimal minimum distances, but be careful, their length $n \leq q$

$\longrightarrow$ There are sharper bounds when $q$ is fixed and $n$ grows!

**Hamming bound:**

For any code $\mathcal{C} \subseteq \mathbb{F}_q^n$,

$$\sharp\mathcal{C} \cdot \sharp\mathcal{B}\left(\mathbf{0}, \left\lfloor \frac{d_{\min}(\mathcal{C})-1}{2} \right\rfloor\right) = \sharp\mathcal{C} \cdot \left(\sum_{r=0}^{\left\lfloor \frac{d_{\min}(\mathcal{C})-1}{2} \right\rfloor} \binom{n}{r}(q-1)^r\right) \leq q^n$$

It asymptotic form when $n \to +\infty$: for any sequence of codes $\mathcal{C}_n \subseteq \mathbb{F}_q^n$ such that the following limits exist:

$$\delta \stackrel{\text{def}}{=} \lim_{n \to \infty} \frac{d_{\min}(\mathcal{C})}{n} \quad \text{and} \quad R \stackrel{\text{def}}{=} \lim_{n \to +\infty} \frac{\log_q \sharp\mathcal{C}_n}{n}$$

we have,

$$\frac{\delta}{2} \leq h_q^{-1}(1-R) \quad \text{where } h_q(x) \stackrel{\text{def}}{=} -(1-x)\log_q(1-x) - x\log_q \frac{x}{q-1}$$

**Proof:**

The Hamming bound is a consequence of the fact that balls centered at codewords with radius $\left\lfloor \frac{d_{\min}(\mathcal{C})-1}{2} \right\rfloor$ are disjoint

The asymptotic form comes from the fact that for a fixed $q$, $\binom{n}{r}(q-1)^n = \text{poly}(n) \cdot q^{nh_q(r/n)}$

**Gilbert-Varshamov bound:**

$q^k \cdot \sharp \mathcal{B}(0, d-2) = q^k \cdot \sum_{i=0}^{d-2} \binom{n}{i} (q-1)^i < q^n \implies$ it exits an $[n, k]_q$-code with minimum distance $d$

$\longrightarrow$ The maximum $d$ reaching the inequality is $d_{\text{GV}}(n, k)$

**Be careful:**

The Gilbert-Varshamov bound states that it exists an $[n, k]_q$-code $\mathcal{C}$ with $d_{\min}(\mathcal{C}) \geq d_{\text{GV}}(n, k)$, not that for all $[n, k]_q$-code $\mathcal{C}$, $d_{\min}(\mathcal{C}) \leq d_{\text{GV}}(n, k)$

*Almost all codes reach **asymptotically** the Gilbert-Varshamov bound*

**Asymptotic Gilbert-Varshamov bound:**

Let $\varepsilon > 0$ and $\delta_{GV} = h_q^{-1}(1 - R)$. We have for uniform $[n, Rn]_q$-codes $\mathcal{C}$,

$$\mathbb{P}_\mathcal{C}\left((1 - \varepsilon)\delta_{GV} < \frac{d_{\min}(\mathcal{C})}{n} < (1 + \varepsilon)\delta_{GV}\right) \geq 1 - q^{-\alpha n(1 + o(1))}$$

where $\alpha \stackrel{\text{def}}{=} \min\left((1 - R) - h_q\left((1 + \varepsilon)\delta_{GV}\right), h_q\left((1 - \varepsilon)\delta_{GV}\right) - (1 - R)\right) > 0$
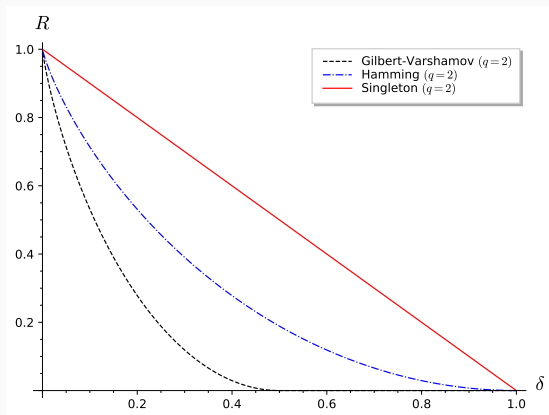
**1M\$ open question:**

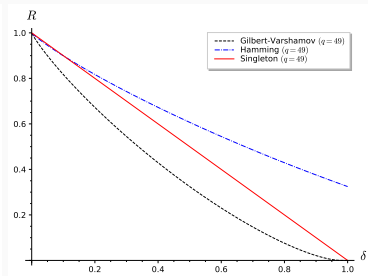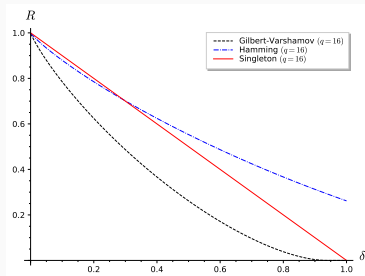Does it exist a sequence of binary linear codes $\mathcal{C}_n$ with rate $R$ such that

$$\frac{d_{\min}(\mathcal{C}_n)}{n} \xrightarrow[n \to +\infty]{} \delta > \delta_{GV} = h_2^{-1}(1 - R)?$$

43

Bounds for sequences of $[n, Rn]_q$-codes $\mathcal{C}_n$ s.t. $\frac{d_{\min}(\mathcal{C}_n)}{n} \xrightarrow[n \to +\infty]{} \delta$ $\left(\text{but } q \text{ is fixed}\right)$

$\left(\text{It exists codes above Gilbert-Varshamov, all codes are below Hamming and Singleton}\right)$

Bounds for sequences of $[n, Rn]_q$-codes $\mathcal{C}_n$ s.t. $\dfrac{d_{\min}(\mathcal{C}_n)}{n} \xrightarrow[n \to +\infty]{} \delta$ $\left(\text{but } q \text{ is fixed}\right)$

We have seen that Reed-Solomon codes reach the Singleton bound $\left(red\ curve\right)$

But the Hamming bound $\left(blue\ curve\right)$ is an upper-bound below the Singleton bound

Don't forget that for Reed-Solomon codes we have $q \geq n$ and for our curve we let $n \rightarrow +\infty$ while $q$ is fixed!

# DECODING REED-SOLOMON CODES

**Reed-Solomon (RS) codes:**

$\mathbf{x} \in \mathbb{F}_q^n$ such that $x_i \neq x_j$ $\left(\text{in particular } n \leq q\right)$ and $k \leq n$. The code $RS_k(\mathbf{x})$ is defined as

$$RS_k(\mathbf{x}) \stackrel{\text{def}}{=} \left\{(f(x_1), \ldots, f(x_n)) : f \in \mathbb{F}_q[X] \text{ and } \deg(f) < k\right\}$$

$\longrightarrow$ These codes are used in QR-codes!

**Exercise:**

Show that $RS_k(\mathbf{x})$ has generator matrix

$$\mathbf{G} \stackrel{\text{def}}{=} \begin{pmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_n \\ x_1^2 & x_2^2 & \cdots & x_n^2 \\ \vdots & \vdots & \vdots & \vdots \\ x_1^{k-1} & x_2^{k-1} & \cdots & x_n^{k-1} \end{pmatrix}$$

Decoding algorithm:

Given, $RS_k(x)$ and $c + e$ such that $\begin{cases} c \in RS_k(x) \\ |e| \leq \left\lfloor \frac{n-k}{2} \right\rfloor \end{cases}$

Then, we can efficiently recover $(c, e)$

Given $\mathbf{y} = \mathbf{c} + \mathbf{e}$ where $\begin{cases} \mathbf{c} \in \mathrm{RS}_k(\mathbf{x}) \\ |\mathbf{e}| \leq \left\lfloor \frac{n-k}{2} \right\rfloor \end{cases}$ .

By definition, $\mathbf{c} = \left( f(x_i) \right)_i$ where $f \in \mathbb{F}_q[X]$ is unknown with $\deg(f) < k$

1. Let $\mathcal{I}$ be the unknown set of positions where $e_i \neq 0$, i.e.,

$$\mathcal{I} = \left\{ i \in [1, n] : \ y_i \neq f(x_i) \right\}$$

**Fundamental idea (I):**

Let $E \in \mathbb{F}_q[X]$ be the following unknown polynomial,

$$E(X) = \prod_{i \in \mathcal{I}} (X - x_i) \ \text{ which has degree} \leq \left\lfloor \frac{n-k}{2} \right\rfloor \text{ by assumption on } |\mathbf{e}|$$

2. By definition of $\mathcal{I}$ and $E$,

$$\forall i \in [1, n], \ y_i E(x_i) = f(x_i) E(x_i)$$

3. The $x_i$'s and $y_i$'s are known: we have above a quadratic system to solve which is a priori not easy

2. By definition of $\mathcal{I}$ and $E$,

$$\forall i \in [1, n], \quad y_i E(x_i) = f(x_i) E(x_i)$$

3. The $x_i$'s and $y_i$'s are known, we have above a quadratic system to solve which is not easy

**Fundamental idea (II): linearize**

Solve the following linear system for unknown $N \in \mathbb{F}_q[X]$ with degree $\leq k - 1 + \left\lfloor \frac{n-k}{2} \right\rfloor$,

$$\forall i \in [1, n], \quad y_i E(x_i) = N(x_i)$$

There are $n$ equations and $k + 2 \left\lfloor \frac{n-k}{2} \right\rfloor + 1$ unknowns $\Big($ coefficients of $N$ and $E\Big)$

$\longrightarrow (E, Ef)$ is a solution but it it is not the only one. . .

### Fundamental idea (II): linearize

Solve the following linear system for **unknown** $N \in \mathbb{F}_q[X]$ with degree $\leq k - 1 + \left\lfloor \frac{n-k}{2} \right\rfloor$ and $E \in \mathbb{F}_q[X]$ with degree $\leq \left\lfloor \frac{n-k}{2} \right\rfloor$,

$$\forall i \in [1, n], \ y_i E(x_i) = N(x_i) \tag{1}$$

### Lemma:

Any non-zero solution $(E_1, N_1)$ and $(E_2, N_2)$ of the above system is such that $\frac{N_1}{E_1} = \frac{N_2}{E_2} = f$

$\longrightarrow$ Therefore the decoding algorithm only consists in computing an non-zero solution $(E, N)$ and to output $f = N/E$

### Proof:

First, if $E_i = 0$, then by Equation (1), $N_i$ has $n > k - 1 + \lfloor (n - k)/2 \rfloor$ zeros and $N_i = 0$. Therefore, $E_i \neq 0$. Now set $R = N_1 E_2 - N_2 E_1$. We have,

$$\mathbf{deg}(R) \leq k - 1 + 2 \left\lfloor \frac{n-k}{2} \right\rfloor \leq n - 1$$

On the other hand, by Equation (1),

$$\forall i \in [1, n], \ R(x_i) = N_1(x_i) E_2(x_i) - N_2(x_i) E_1(x_i) = y_i E_1(x_i) E_2(x_i) - y_i E_1(x_i) E_2(x_i) = 0$$

Therefore, $R = 0$, showing $N_1/E_1 = N_2/E_2$. But $(N, Ef)$ is a non-zero solution concluding the proof

We have demonstrated that we can decode Reed-Solomon codes. Does it exist other codes that we know how to decode?

$$\longrightarrow \text{Yes!}$$

▶ **Algebraic decoders:** Reed-Solomon, alternant, geometric codes

▶ **Probabilistic decoders:** LDPC, Turbo, Polar codes

**Be careful:**
- It is an hard problem to design codes with an efficient decoding algorithm
- When designing a decoding algorithm we have to be cautious about the parameters: the field size $q$ or the decoding distance $\Big($larger it is, harder is to decode$\Big)$

▶ An introduction to Low-Density Parity-Check $\left(\text{LDPC}\right)$ codes available here:

https://repository.arizona.edu/handle/10150/607470

EXERCISE SESSION