

## Post-Quantum Cryptography - Codes

### Lecture 2: Random Codes

Lecturer: Thomas Debris-Alazard

#### INTRODUCTION

We study in this course *random codes*, *i.e.* codes whose parity-check or generator matrix is drawn uniformly at random. However, in light of the history of error correcting codes which has consisted in finding codes becoming more and more complex and structured, one may wonder but why then are we wasting our time to study random codes? That may come as a surprise but random codes enlighten about what we could expect or not in a simple fashion, and even better what is optimal or not. The most prominent example of the interest of random codes is the famous Shannon theorem about the capacity of some “noisy channels”. Roughly speaking, Shannon gave (for some error models) the maximum amount of errors that “can” be theoretically decoded with codes of fixed rate. Shannon made his proof by using random codes and he has shown that they are precisely those which reach the optimality.

Our aim in these lecture notes is to study carefully these kind of codes and to show that they enable to answer many questions like for instance:

- How many vectors of Hamming weight  $t$  do we expect in a code?
- What is the typical minimum distance of a code?
- etc...

Our study will have an important consequence for cryptographic purposes: a better understanding of the Decoding Problem  $\text{DP}(n, q, R, \tau)$  that was defined in lecture notes 1. We will be able to predict with a very good accuracy the number of solutions of this problem as a function of its parameters. This will be particularly helpful to understand the behaviour of algorithms solving it.

#### 1. PREREQUISITES

**Basic notation.** In all these lecture notes,  $q$  will denote a fixed field size while  $R$  will be a *constant* in  $(0, 1)$ . On the other hand,  $\tau(n)$  will denote any function of  $n$  taking its values in  $(0, 1)$ . To simplify notation, since  $n$  is clear from the context, we will drop the dependency in  $n$  and simply write  $\tau$ . Furthermore, parameters  $k$  and  $t$  will always (even implicitly) be defined as  $k \stackrel{\text{def}}{=} \lfloor Rn \rfloor$  and  $t \stackrel{\text{def}}{=} \lfloor \tau n \rfloor$ . A function  $f(n)$  is said to be negligible, and we denote this by  $f \in \text{negl}(n)$ , if for all polynomial  $p(n)$ ,  $|f(n)| < |p(n)|^{-1}$  for all sufficiently large  $n$ .

Many asymptotic results will be given. As all our parameters are functions of  $n$ , our asymptotic results will always hold for:

$$n \longrightarrow +\infty.$$

The parameter  $n$  is in most cryptographic applications roughly given by several thousands.

The following function  $h_q$ , known as the  $q$ -ary entropy, will play an important role:

$$h_q : x \in [0, 1] \longmapsto -x \log_q \left( \frac{x}{q-1} \right) - (1-x) \log_q (1-x) \quad (\text{extended by continuity in } 0 \text{ and } 1).$$

It is equal to the entropy of a random variable  $e$  over  $\mathbb{F}_q$  distributed like the error for a  $q$ -ary symmetric channel of crossover probability  $x$ , *i.e.*  $\mathbb{P}(e=0) = 1-x$  and  $\mathbb{P}(e=\alpha) = \frac{x}{q-1}$  for any  $\alpha \in \mathbb{F}_q^*$ .

It can be verified that  $h_q$  is an increasing function over  $\left[0, \frac{q-1}{q}\right]$  and a decreasing function over  $\left[\frac{q-1}{q}, 1\right]$ . The  $q$ -ary entropy is involved in the estimation of  $\#\mathcal{S}_t$  where  $\mathcal{S}_t$  is defined as the sphere of radius  $t$  for the Hamming distance  $|\cdot|$ , namely

$$\mathcal{S}_t \stackrel{\text{def}}{=} \{\mathbf{x} \in \mathbb{F}_q^n : |\mathbf{x}| = t\}.$$

The following elementary lemma will be at the core of most of our asymptotic results.

**Lemma 1.** *Let  $t \stackrel{\text{def}}{=} \lfloor \tau n \rfloor$ . We have  $\#\mathcal{S}_t = \binom{n}{t}(q-1)^t$  and*

$$q^{n(h_q(\tau) + O(\frac{\log_q(n)}{n}))} \leq \binom{n}{t}(q-1)^t \leq q^{nh_q(\tau)}.$$

*Asymptotically,*

$$\frac{1}{n} \log_q \left( \binom{n}{t}(q-1)^t \right) = h_q(\tau) + O\left(\frac{\log_q n}{n}\right).$$

**Probabilistic notation.** During these lecture notes we wish to emphasize on which probability space the probabilities or the expectations are taken. Therefore we will denote by a subscript the random variable specifying the associated probability space over which the probabilities or expectations are taken. For instance the probability  $\mathbb{P}_X(A)$  of the event  $A$  is taken over  $\Omega$  the probability space over which the random variable  $X$  is defined, *i.e.* if  $X$  is for instance a real random variable,  $X$  is a function from a probability space  $\Omega$  to  $\mathbb{R}$ , and the aforementioned probability is taken according to the measure chosen for  $\Omega$ .

**Statistical distance.** An essential tool for many cryptographic applications is the *statistical distance*, sometimes called the *total variational distance*. It is a distance for probability distributions, which in the case where  $X$  and  $Y$  are two random variables taking their values in a same finite space  $\mathcal{E}$  is defined as

$$(1) \quad \Delta(X, Y) \stackrel{\text{def}}{=} \frac{1}{2} \sum_{a \in \mathcal{E}} |\mathbb{P}(X = a) - \mathbb{P}(Y = a)|.$$

An equivalent definition is given by

$$(2) \quad \Delta(X, Y) \stackrel{\text{def}}{=} \max_{A \subseteq \mathcal{E}} |\mathbb{P}_X(A) - \mathbb{P}_Y(A)|.$$

Depending on the context, (1) or (2) is the most useful. A direct consequence of (2) is that given any event  $A$ , we have  $|\mathbb{P}_X(A) - \mathbb{P}_Y(A)| \leq \Delta(X, Y)$ . Therefore, computing probabilities over  $X$  or  $Y$  will differ by at most  $\Delta(X, Y)$ . Furthermore, given a single observation, coming from  $X$  or  $Y$  with probability  $1/2$ , we will be able to guess which with probability at most  $1/2 + \Delta(X, Y)/2$  and there is a strategy to reach this probability of guessing correctly.

The statistical distance enjoys many interesting properties. Among others, it cannot increase by applying a function  $f$ ,

$$\Delta(f(X), f(Y)) \leq \Delta(X, Y) \quad (\text{data processing inequality}).$$

The function  $f$  can be randomized, but its internal randomness has to be independent from  $X$  and  $Y$  for the data processing inequality to hold. In particular, it implies that the “success” probability of any algorithm  $\mathcal{A}$  for inputs distributed according to  $X$  or  $Y$ , can only differ by at most  $\Delta(X, Y)$ .

In our applications we will focus on distributions  $X, Y$  such that their statistical distance is negligible. It will show (as a consequence of the data processing inequality) that  $X$  and  $Y$  are computationally indistinguishable<sup>(1)</sup> without requiring any computational argument with a reduction.

<sup>(1)</sup>See here for a definition: <https://www.cs.princeton.edu/courses/archive/spr10/cos433/lec4.pdf>

One can define various other distances for capturing in a cryptographic context the differences between two distributions. For instance, we can cite the family Renyi divergences, but this is out of the scope of these lecture notes.

## 2. RANDOM CODES

**The model of random codes.** In these lecture notes we will use two probabilistic models that will be referred to as *random*  $[n, k]_q$ -codes. The first one is by choosing a code  $\mathcal{C}$  by picking uniformly at random a generator matrix  $\mathbf{G} \in \mathbb{F}_q^{k \times n}$  (i.e.  $\mathcal{C} \stackrel{\text{def}}{=} \{\mathbf{m}\mathbf{G} : \mathbf{m} \in \mathbb{F}_q^k\}$ ). However, all the probabilistic results of these lecture notes are easier to prove if, instead, we choose  $\mathcal{C}$  by picking uniformly at random a parity-check matrix  $\mathbf{H} \in \mathbb{F}_q^{(n-k) \times n}$  (i.e.  $\mathcal{C} \stackrel{\text{def}}{=} \{\mathbf{c} \in \mathbb{F}_q^n : \mathbf{H}\mathbf{c}^\top = \mathbf{0}\}$ ). We will denote  $\mathbb{P}_{\mathbf{G}}$  and  $\mathbb{P}_{\mathbf{H}}$  respectively the probabilities in these two models.

It may be pointed out that in both models we don't pick uniformly at random an  $[n, k]_q$ -code. Indeed, the first model always produces codes of dimension  $\leq k$  whereas in the second model codes are always of dimension  $\geq k$ . One may wonder why don't we pick  $\mathbf{G}$  (resp.  $\mathbf{H}$ ) uniformly at random among the  $k \times n$  (resp.  $(n-k) \times n$ ) matrices of rank  $k$  (resp.  $n-k$ )? First, computations are much more complicated in this "exact" model. Furthermore, it turns out that it is pointless. Roughly speaking, the  $\mathbf{G}$ -model produces codes of dimension  $= k$  with probability  $1 - O(q^{-(n-k)})$  while in the  $\mathbf{H}$ -model we get a code of dimension  $= k$  with probability  $1 - O(q^{-k})$ . As shown in the following lemma this result can even be expressed in a stronger way, our probabilistic models are exponentially close, *for the statistical distance*, to the "exact" model. Therefore all our computations in the  $\mathbf{G}$  or  $\mathbf{H}$  models can really be thought as by picking uniformly at random an  $[n, k]_q$ -code  $\mathcal{C}$ .

**Lemma 2.** Let  $\mathbf{G} \in \mathbb{F}_q^{k \times n}$  (resp.  $\mathbf{H} \in \mathbb{F}_q^{(n-k) \times n}$ ) be a uniformly random matrix and  $\mathbf{G}_k \in \mathbb{F}_q^{k \times n}$  (resp.  $\mathbf{H}_{n-k} \in \mathbb{F}_q^{(n-k) \times n}$ ) be a uniformly random matrix of rank  $k$  (resp.  $n-k$ ). We have:

$$\Delta(\mathbf{G}, \mathbf{G}_k) = O\left(q^{-(n-k)}\right) \quad (\text{resp. } \Delta(\mathbf{H}, \mathbf{H}_{n-k}) = O\left(q^{-k}\right)).$$

*Proof.* Let us prove the lemma for  $(\mathbf{G}, \mathbf{G}_k)$ , the other case will be similar. First it is a classical fact that the density of rank  $k$  matrices among  $\mathbb{F}_q^{k \times n}$  is equal to  $1 - O(q^{-(n-k)})$ . Therefore, given some rank  $k$  matrix  $\mathbf{R} \in \mathbb{F}_q^{k \times n}$ , we have:

$$\mathbb{P}(\mathbf{G}_k = \mathbf{R}) = \frac{1}{q^{k \times n} (1 - O(q^{-(n-k)}))}.$$

It leads to the following computation:

$$\begin{aligned} 2\Delta(\mathbf{G}, \mathbf{G}_k) &= \sum_{\substack{\mathbf{R} \in \mathbb{F}_q^{k \times n} \\ \text{rank}(\mathbf{R})=k}} |\mathbb{P}(\mathbf{G} = \mathbf{R}) - \mathbb{P}(\mathbf{G}_k = \mathbf{R})| + \sum_{\substack{\mathbf{R} \in \mathbb{F}_q^{k \times n} \\ \text{rank}(\mathbf{R}) \neq k}} \mathbb{P}(\mathbf{G} = \mathbf{R}) \\ &= \sum_{\substack{\mathbf{R} \in \mathbb{F}_q^{k \times n} \\ \text{rank}(\mathbf{R})=k}} \left| \frac{1}{q^{k \times n}} \left( 1 - \frac{1}{1 - O(q^{-(n-k)})} \right) \right| + \sum_{\substack{\mathbf{R} \in \mathbb{F}_q^{k \times n} \\ \text{rank}(\mathbf{R}) \neq k}} \frac{1}{q^{k \times n}} \\ &= O\left(q^{-(n-k)}\right) \end{aligned}$$

which concludes the proof.  $\square$

Now one may wonder why do we consider two models for random  $[n, k]_q$ -codes? It turns out that depending of the context, computations might be easier and/or more natural in one model rather than in the other one. In addition, for the same reasons as those given in the previous lemma,  $\mathbf{G}$  and  $\mathbf{H}$  models are closely related,

computations in both probabilistic models will outcome the same results up to an additive exponentially small factor.

**Lemma 3.** *Let  $\mathcal{E}$  be a set of linear codes of length  $n$  in  $\mathbb{F}_q$  which is defined as an event. We have,*

$$|\mathbb{P}_{\mathbf{G}}(\mathcal{E}) - \mathbb{P}_{\mathbf{H}}(\mathcal{E})| = O\left(q^{-\min(k, n-k)}\right).$$

*Proof.* Let  $\mathbf{G}_k$  and  $\mathbf{H}_{n-k}$  be defined as in Lemma 2. Notice that  $\mathbb{P}_{\mathbf{G}_k}(\mathcal{E}) = \mathbb{P}_{\mathbf{H}_{n-k}}(\mathcal{E})$ , in both models, we exactly pick uniformly at random an  $[n, k]_q$ -code. It leads to the following computation:

$$\begin{aligned} |\mathbb{P}_{\mathbf{G}}(\mathcal{E}) - \mathbb{P}_{\mathbf{H}}(\mathcal{E})| &\leq |\mathbb{P}_{\mathbf{G}}(\mathcal{E}) - \mathbb{P}_{\mathbf{G}_k}(\mathcal{E})| + |\mathbb{P}_{\mathbf{H}_{n-k}}(\mathcal{E}) - \mathbb{P}_{\mathbf{H}}(\mathcal{E})| \\ &\leq \Delta(\mathbf{G}, \mathbf{G}_k) + \Delta(\mathbf{H}, \mathbf{H}_{n-k}) \end{aligned}$$

where in the last line we used Equation (2). It concludes the proof by using Lemma 2.  $\square$

**Exercise 1.** *Let us introduce the following variant of DP (which has been introduced in Lecture 1) with generator matrices instead of parity-check matrices*

$\text{DP}'(n, q, R, \tau)$ . Let  $k \stackrel{\text{def}}{=} \lfloor Rn \rfloor$  and  $t \stackrel{\text{def}}{=} \lfloor \tau n \rfloor$ .

- Input:  $(\mathbf{G}, \mathbf{y} \stackrel{\text{def}}{=} \mathbf{s}\mathbf{G} + \mathbf{x})$  where  $\mathbf{G}, \mathbf{s}$  and  $\mathbf{x}$  are uniformly distributed over  $\mathbb{F}_q^{k \times n}$ ,  $\mathbb{F}_q^k$  and words of Hamming weight  $t$  in  $\mathbb{F}_q^n$ .
- Output: an error  $\mathbf{e} \in \mathbb{F}_q^n$  of Hamming weight  $t$  such that  $\mathbf{y} - \mathbf{e} = \mathbf{m}\mathbf{G}$  for some  $\mathbf{m} \in \mathbb{F}_q^k$ .

Show that for any algorithm  $\mathcal{A}$  solving this problem with probability  $\varepsilon$  and time  $T$ , there exists an algorithm  $\mathcal{B}$  which solves  $\text{DP}(n, q, R, \tau)$  in time  $O(n^3 + T)$  with probability  $\geq \varepsilon - O(q^{-\min(k, n-k)})$ . Show that we can exchange  $\text{DP}'$  by  $\text{DP}$  in the previous question.

**Remark 1.** *The above exercise shows that defining DP with generator or parity-check matrices is just a matter of personal taste, it does not change the average hardness.*

**A first computation with random codes.** Now that random codes are well defined, we are ready to make our first computation in this probabilistic model. The following elementary lemma gives the probability (over the codes) that a fixed non-zero word  $\mathbf{y}$  reaches some syndrome  $\mathbf{s}$  according to the code. In particular, by setting  $\mathbf{s}$  to  $\mathbf{0}$ , we obtain the probability that  $\mathbf{y}$  belongs to the code. This lemma will be at the core of all our results about random codes.

**Lemma 4.** *Given  $\mathbf{s} \in \mathbb{F}_q^{n-k}$  and  $\mathbf{y} \in \mathbb{F}_q^n$  such that  $\mathbf{y} \neq \mathbf{0}$ , we have for  $\mathbf{H}$  being uniformly distributed at random in  $\mathbb{F}_q^{(n-k) \times n}$ ,*

$$\mathbb{P}_{\mathbf{H}}(\mathbf{y}\mathbf{H}^T = \mathbf{s}) = \frac{1}{q^{n-k}}.$$

*Proof.* Let  $h_{i,j}$  be the coefficient of  $\mathbf{H}$  at position  $(i, j)$ . Without loss of generality, we can suppose that  $y_1 = 1$  (by permuting  $\mathbf{y}$  if  $y_1 = 0$  and then multiplying by  $y_1^{-1}$  which is possible as we work in  $\mathbb{F}_q$ ). The probability that we are looking for is the probability of the following event:

$$\forall i \in \llbracket 1, n-k \rrbracket, \quad h_{i,1} = s_i - \sum_{j=2}^n h_{i,j}y_j$$

Recall that  $\mathbf{H}$  is uniformly distributed: the  $h_{i,j}$ 's are independent and equidistributed. Therefore the above  $n-k$  equations will be independently true with probability  $1/q$  which concludes the proof.  $\square$

**Exercise 2.** *Show that for any non-zero  $\mathbf{y} \in \mathbb{F}_q^n$ ,*

$$\mathbb{P}_{\mathbf{G}}(\mathbf{y} \in \mathbb{C}^*) = \frac{1}{q^k}.$$

### 3. WEIGHT DISTRIBUTION OF COSETS OF RANDOM CODES

The aim of this section is to answer the following question: given a random code  $\mathcal{C}$  and a fixed vector  $\mathbf{y} \in \mathbb{F}_q^n$ , how many codewords  $\mathbf{c} \in \mathcal{C}$  do we expect to be at Hamming distance  $t$  from  $\mathbf{y}$ ? Or equivalently, given a parity-check matrix of our random code  $\mathcal{C}$  and a fixed syndrome  $\mathbf{s}$ , how many vectors  $\mathbf{e}$  of Hamming weight  $t$  do we expect to reach the syndrome  $\mathbf{s}$  according to  $\mathbf{H}$ ? Notice that deriving an answer to these questions in the particular cases  $\mathbf{y} = \mathbf{0}$  and  $\mathbf{s} = \mathbf{0}$  enables to compute the expected number of codewords of weight  $t$  in  $\mathcal{C}$ . It will be useful to compute the expected minimum distance of a code.

These results will have an important consequence: a better understanding of the Decoding Problem (DP) that we recall now.

**Problem 1** (Decoding Problem -  $\text{DP}(n, q, R, \tau)$ ). Let  $k \stackrel{\text{def}}{=} \lfloor Rn \rfloor$  and  $t \stackrel{\text{def}}{=} \lfloor \tau n \rfloor$ .

- Input :  $(\mathbf{H}, \mathbf{s} \stackrel{\text{def}}{=} \mathbf{xH}^\top)$  where  $\mathbf{H}$  (resp.  $\mathbf{x}$ ) is uniformly distributed over  $\mathbb{F}_q^{(n-k) \times n}$  (resp.  $\mathcal{S}_t$ , words of Hamming weight  $t$  in  $\mathbb{F}_q^n$ ).
- Output : an error  $\mathbf{e} \in \mathbb{F}_q^n$  of Hamming weight  $t$  such that  $\mathbf{eH}^\top = \mathbf{s}$ .

According to our probabilistic model, this problem really corresponds to decode a random  $[n, k]_q$ -code of parity-check matrix  $\mathbf{H}$ . In that case it is natural to wonder how many vectors  $\mathbf{e} \in \mathcal{S}_t$  are expected to reach the syndrome  $\mathbf{s}$  according to  $\mathbf{H}$ , but why? To understand this let us take a toy example. A trivial solution to solve DP is to pick a random error  $\mathbf{e} \in \mathcal{S}_t$  with the hope that it gives a solution. By definition there is a solution to our problem (here  $\mathbf{x}$ ). If there is exactly one solution, our success probability is given by  $\frac{1}{\binom{n}{t}(q-1)^t}$ . But now imagine that we expect  $N$  solutions to our problem. In that case we would expect our success probability to be equal to  $\approx \frac{N}{\binom{n}{t}(q-1)^t}$ . It is therefore important to know the value of  $N$  to be able to predict the running time of our algorithm. It is the aim of what follows.

**Notation.** Given  $\mathbf{H} \in \mathbb{F}_q^{(n-k) \times n}$  and  $\mathbf{s} \in \mathbb{F}_q^{n-k}$ , let

$$N_t(\mathcal{C}, \mathbf{s}) \stackrel{\text{def}}{=} \#\{\mathbf{e} \in \mathcal{S}_t : \mathbf{eH}^\top = \mathbf{s}\}$$

where implicitly  $\mathcal{C}$  is defined as  $\{\mathbf{c} \in \mathbb{F}_q^n : \mathbf{Hc}^\top = \mathbf{0}\}$ . Notice that  $N_t(\mathcal{C}, \mathbf{s})$  is a random variable that gives the number of solutions of DP with input  $(\mathbf{H}, \mathbf{s})$  (where  $\mathbf{s} \in \mathbb{F}_q^{n-k}$  is fixed and not necessarily computed as some  $\mathbf{xH}^\top$  for  $\mathbf{x} \in \mathcal{S}_t$ ). On the other hand,  $N_t(\mathcal{C}, \mathbf{0})$  is a random variable that gives the number of codewords  $\mathbf{c} \in \mathcal{C}$  of Hamming weight  $t$ .

**Expected weight distribution of cosets.** From now on, our main objective is to compute the expected number of Hamming weight  $t$  vectors in a given coset, namely to compute the expectation of  $N_t(\mathcal{C}, \mathbf{s})$  over  $\mathcal{C}$ . To avoid any suspense, we will prove that for any syndrome  $\mathbf{s} \in \mathbb{F}_q^{n-k}$ , the expectation of  $N_t(\mathcal{C}, \mathbf{s})$  is given by  $\binom{n}{t}(q-1)^t/q^{n-k}$ . However, before showing this result, let us start to understand how this quantity behaves as function of its parameters, namely  $\tau = t/n$  and  $R = k/n$ . By Lemma 1, we have

$$(3) \quad \frac{1}{n} \log_q \binom{n}{t} (q-1)^t / q^{n-k} = h_q(\tau) - (1-R) + O\left(\frac{\log_q n}{n}\right).$$

Recall now that  $x \in [0, 1] \mapsto h_q(x)$  is an increasing function over  $\left[0, \frac{q-1}{q}\right]$  and a decreasing function over  $\left[\frac{q-1}{q}, 1\right]$ . Furthermore,  $h_q(0) = 0$  and  $h_q(1) = \log_q(q-1)$ . This shows that  $N_t(\mathcal{C}, \mathbf{s})$  is expected (according to  $\tau$ ) to be exponentially small or large (in  $n$ ) at the exception of one value  $\tau^-$  and potentially a second one in the case where  $(1-R) \geq \log_q(q-1)$ , that we will denote  $\tau^+$ . We summarize the picture by drawing in Figure 1 the logarithm in basis 3 (for  $n$  large enough) of  $\binom{n}{t}(q-1)^t/q^{n-k}$  when  $q = 3$  and  $k/n = 1/4$ .

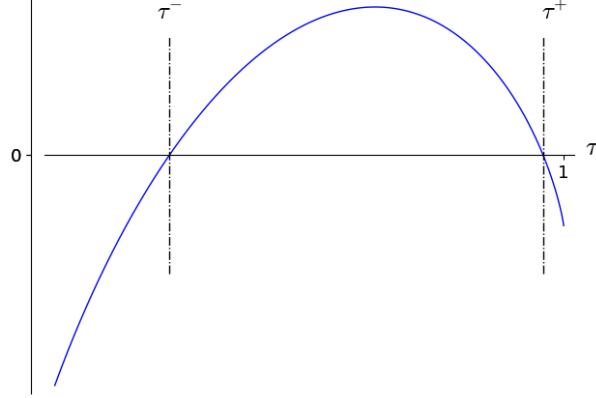


FIGURE 1.  $\lim_{n \rightarrow +\infty} \frac{1}{n} \log_q \binom{n}{t} (q-1)^t / q^{n-k}$  when  $q = 3$  and  $k/n = 1/4$  as function of  $\tau = t/n$ .

It turns out that an analytic expression of  $\tau^-$  and  $\tau^+$  can be given,

$$(4) \quad \tau^- \stackrel{\text{def}}{=} g_q^-(1-R) \quad \text{and} \quad \tau^+ \stackrel{\text{def}}{=} g_q^+(1-R) \quad \text{when } R \leq 1 - \log_q(q-1)$$

where  $g_q^-$  (resp.  $g_q^+$ ) denotes the inverse of  $h_q$  over  $\left[0, \frac{q-1}{q}\right]$  (resp.  $\left[\frac{q-1}{q}, 1\right]$ ).

**Remark 2.** As we will see in Section 4,  $\tau^-$  is commonly called the relative Gilbert-Varshamov distance or bound.

Quantities  $\tau^-$  and  $\tau^+$  give the boundaries between which we expect  $N_t(\mathbb{C}, \mathbf{s})$  to be exponentially large as we show now.

**Proposition 1.** Let  $k \stackrel{\text{def}}{=} \lfloor Rn \rfloor$ ,  $t \stackrel{\text{def}}{=} \lfloor \tau n \rfloor$  and  $\mathbf{s} \in \mathbb{F}_q^{n-k}$ . We have:

$$(5) \quad \mathbb{E}_{\mathbf{H}}(N_t(\mathbb{C}, \mathbf{s})) = \frac{\binom{n}{t} (q-1)^t}{q^{n-k}}.$$

When  $\tau \in \begin{cases} (\tau^-, \tau^+) & \text{if } R \leq 1 - \log_q(q-1) \\ (\tau^-, 1) & \text{otherwise} \end{cases}$ , we expect  $N_t(\mathbb{C}, \mathbf{s})$  to be exponentially large:

$$\mathbb{E}_{\mathbf{H}}(N_t(\mathbb{C}, \mathbf{s})) = q^{\alpha n(1+o(1))} \quad \text{where} \quad \alpha \stackrel{\text{def}}{=} h_q(\tau) - 1 + R > 0.$$

In the case where  $\tau = \begin{cases} \tau^- & \text{or } \tau^+ \\ \tau^- & \text{otherwise.} \end{cases}$  if  $R \leq 1 - \log_q(q-1)$ , the expectation of  $N_t(\mathbb{C}, \mathbf{s})$  equals  $P(n)(1+o(1))$  for some polynomial  $P$ .

*Proof.* Let  $\mathbf{1}_{\mathbf{e}}$  be the indicator function of the event “ $\mathbf{eH}^\top = \mathbf{s}$ ”. It is readily verified that by definition,

$$N_t(\mathbb{C}, \mathbf{s}) = \sum_{\mathbf{e} \in \mathcal{S}_t} \mathbf{1}_{\mathbf{e}}.$$

We have the following computation,

$$\begin{aligned}
\mathbb{E}_{\mathbf{H}}(N_t(\mathcal{C}, \mathbf{s})) &= \mathbb{E}_{\mathbf{H}} \left( \sum_{\mathbf{e} \in \mathcal{S}_t} \mathbf{1}_{\mathbf{e}} \right) \\
&= \sum_{\mathbf{e} \in \mathcal{S}_t} \mathbb{E}_{\mathbf{H}}(\mathbf{1}_{\mathbf{e}}) \quad (\text{by linearity of the expectation}) \\
&= \sum_{\mathbf{e} \in \mathcal{S}_t} \mathbb{P}_{\mathbf{H}}(\mathbf{eH}^T = \mathbf{s})
\end{aligned}$$

which gives (5) by using Lemma 4. The second part of the proposition is a consequence of Lemma 1.  $\square$

**Exercise 3.** Show that the average number of solutions of  $\text{DP}(n, q, R, \tau)$  (over the input distribution) is given by

$$1 + \frac{\binom{n}{t}(q-1)^t - 1}{q^{n-k}}$$

**Remark 3.** Proposition 1 can actually be stated more generally. Given some set  $\mathcal{E} \subseteq \mathbb{F}_q^n$ , we can show that the expected number of vectors  $\mathbf{e} \in \mathcal{E}$  that reach some syndrome for a random  $[n, k]_q$ -code is given by  $\#\mathcal{E}/q^{n-k}$ .

**Remark 4.** The expected number of codewords of weight  $t$  in a random  $[n, k]_q$ -code is given by  $\binom{n}{t}(q-1)^t/q^{n-k}$  (by setting  $\mathbf{s}$  to  $\mathbf{0}$  in Proposition 1). This statement, in the same manner as in the previous remark, can be generalized to give the expected number of codewords in any set  $\mathcal{E} \subseteq \mathbb{F}_q^n$ . In particular, it can be used to obtain the expected number of codewords of “weight”  $t$  that belong to a random code, for any notion of weight and therefore any metric.

**Exercise 4.** Show that ( $t > 0$ ),

$$\mathbb{E}_{\mathbf{G}}(\#\{\mathbf{m} \in \mathbb{F}_q^k : |\mathbf{mG}| = t\}) = \frac{q^k - 1}{q^n} \binom{n}{t} (q-1)^t \quad \text{and} \quad \mathbb{E}_{\mathbf{H}}(\#\{\mathbf{c} \in \mathcal{C} : |\mathbf{c}| \text{ is odd}\}) = \frac{1}{2} \frac{q^n - (2-q)^n}{q^{n-k}}.$$

**Hint:** For the first part of the exercise first show that  $\mathbf{mG}$  is uniformly distributed over  $\mathbb{F}_q^n$  when  $\mathbf{m} \in \mathbb{F}_q^k \setminus \{\mathbf{0}\}$ .

At this point our work has given the *expected* number of solutions of  $\text{DP}(n, q, R, \tau)$ . Situation is depicted in Figure 2.

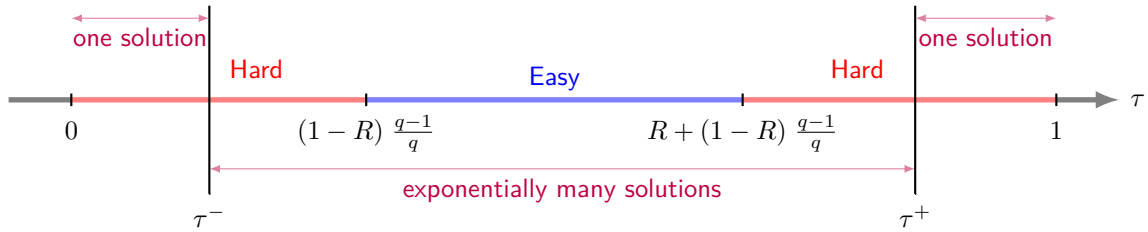


FIGURE 2. Hardness and expected number of solutions of  $\text{DP}(n, q, R, \tau)$  as function of  $\tau$ .

However, can we be much more precise? For instance, can we give with an overwhelming probability, and therefore for almost all codes, the number of solutions of  $\text{DP}$ ? The answer is yes. Below is given two techniques for achieving this, the first one uses Markov's inequality (*first moment technique*) and the second one, which is more accurate, uses Bienaymé-Tchebychev's inequality (*second moment technique*).

**Proposition 2** (First Moment Technique). *Let  $\mathbf{s} \in \mathbb{F}_q^{n-k}$ . For any  $a > 0$ , we have*

$$\mathbb{P}_{\mathbf{H}}(N_t(\mathcal{C}, \mathbf{s}) > a) \leq \frac{1}{a} \frac{\binom{n}{t}(q-1)^t}{q^{n-k}}.$$

*Proof.* By Proposition 1,  $\mathbb{E}_{\mathbf{H}}(N_t(\mathcal{C}, \mathbf{s})) = \frac{\binom{n}{t}(q-1)^t}{q^{n-k}}$ . It remains to apply Markov's inequality to conclude the proof.  $\square$

Notice that Proposition 2 is not very accurate. One has to choose  $a \gg \frac{\binom{n}{t}(q-1)^t}{q^{n-k}}$  to obtain a negligible probability that  $N_t(\mathcal{C}, \mathbf{s})$  is larger than  $a$ . But the expectation of  $N_t(\mathcal{C}, \mathbf{s})$  is exactly given by  $\frac{\binom{n}{t}(q-1)^t}{q^{n-k}}$ . A meaningful result would be

$$\mathbb{P}_{\mathbf{H}} \left( \left| N_t(\mathcal{C}, \mathbf{s}) - \frac{\binom{n}{t}(q-1)^t}{q^{n-k}} \right| > a \right) < \varepsilon,$$

for some  $\varepsilon \in \text{negl}(n)$  and  $a$  be such that  $a \in \frac{\binom{n}{t}(q-1)^t}{q^{n-k}} \text{negl}(n)$ . It is precisely the aim of the following proposition.

**Proposition 3** (Second Moment Technique). *Let  $\mathbf{s} \in \mathbb{F}_q^{n-k}$ . For any  $a > 0$ , we have*

$$\mathbb{P}_{\mathbf{H}} \left( \left| N_t(\mathcal{C}, \mathbf{s}) - \frac{\binom{n}{t}(q-1)^t}{q^{n-k}} \right| \geq a \right) \leq \frac{(q-1)\binom{n}{t}(q-1)^t}{a^2 q^{n-k}}.$$

*Proof.* Let  $\mathbf{1}_{\mathbf{e}}$  be the indicator function of the event “ $\mathbf{eH}^T = \mathbf{s}$ ”. By using Bienaymé-Tchebychevs inequality with the random variable  $N_t(\mathcal{C}, \mathbf{s}) = \sum_{\mathbf{e} \in \mathcal{S}_t} \mathbf{1}_{\mathbf{e}}$ , we obtain

$$\begin{aligned} \mathbb{P}_{\mathbf{H}} \left( \left| N_t(\mathcal{C}, \mathbf{s}) - \frac{\binom{n}{t}(q-1)^t}{q^{n-k}} \right| \geq a \right) &\leq \frac{\text{Var}_{\mathbf{H}}(N_t(\mathcal{C}, \mathbf{s}))}{a^2} \\ &= \frac{1}{a^2} \left( \sum_{\mathbf{e} \in \mathcal{S}_t} \text{Var}_{\mathbf{H}}(\mathbf{1}_{\mathbf{e}}) + \sum_{\substack{\mathbf{x}, \mathbf{y} \in \mathcal{S}_t \\ \mathbf{x} \neq \mathbf{y}}} \mathbb{E}_{\mathbf{H}}(\mathbf{1}_{\mathbf{x}} \mathbf{1}_{\mathbf{y}}) - \mathbb{E}_{\mathbf{H}}(\mathbf{1}_{\mathbf{x}}) \mathbb{E}_{\mathbf{H}}(\mathbf{1}_{\mathbf{y}}) \right) \\ &\leq \frac{1}{a^2} \left( \sum_{\mathbf{e} \in \mathcal{S}_t} \mathbb{E}_{\mathbf{H}}(\mathbf{1}_{\mathbf{e}}) + \sum_{\substack{\mathbf{x}, \mathbf{y} \in \mathcal{S}_t \\ \mathbf{x} \neq \mathbf{y}}} \mathbb{E}_{\mathbf{H}}(\mathbf{1}_{\mathbf{x}} \mathbf{1}_{\mathbf{y}}) - \mathbb{E}_{\mathbf{H}}(\mathbf{1}_{\mathbf{x}}) \mathbb{E}_{\mathbf{H}}(\mathbf{1}_{\mathbf{y}}) \right) \\ (6) \quad &= \frac{1}{a^2} \left( \frac{\binom{n}{t}(q-1)^t}{q^{n-k}} + \sum_{\substack{\mathbf{x}, \mathbf{y} \in \mathcal{S}_t \\ \mathbf{x} \neq \mathbf{y}}} \mathbb{E}_{\mathbf{H}}(\mathbf{1}_{\mathbf{x}} \mathbf{1}_{\mathbf{y}}) - \mathbb{E}_{\mathbf{H}}(\mathbf{1}_{\mathbf{x}}) \mathbb{E}_{\mathbf{H}}(\mathbf{1}_{\mathbf{y}}) \right) \end{aligned}$$

where we used that  $\text{Var}_{\mathbf{H}}(\mathbf{1}_{\mathbf{e}}) \leq \mathbb{E}_{\mathbf{H}}(\mathbf{1}_{\mathbf{e}}^2) = \mathbb{E}_{\mathbf{H}}(\mathbf{1}_{\mathbf{e}})$ . Let us now upper-bound the second term of the inequality. To this aim let us prove the following lemma.

**Lemma 5.** *We have*

$$\mathbb{E}_{\mathbf{H}}(\mathbf{1}_{\mathbf{x}} \mathbf{1}_{\mathbf{y}}) \leq \begin{cases} 1/q^{n-k} & \text{if } \mathbf{x} \text{ and } \mathbf{y} \text{ are colinear} \\ 1/q^{2(n-k)} & \text{otherwise.} \end{cases}$$

*Proof.* The result is clear when  $\mathbf{x}$  and  $\mathbf{y}$  are colinear by Lemma 4. Let us suppose that  $\mathbf{x}$  and  $\mathbf{y}$  are not colinear and define

$$\varphi : \mathbf{h} \in \mathbb{F}_q^n \longmapsto (\mathbf{h} \cdot \mathbf{x}, \mathbf{h} \cdot \mathbf{y}).$$



It is readily verified that this linear application has a kernel of  $\mathbb{F}_q$ -dimension  $n - 2$ . Therefore, for any  $a, b \in \mathbb{F}_q$  and  $\mathbf{h} \in \mathbb{F}_q^n$  being uniformly distributed,

$$(7) \quad \mathbb{P}_{\mathbf{h}}(\varphi(\mathbf{h}) = (a, b)) = \frac{q^{n-2}}{q^n} = \frac{1}{q^2}$$

Let us remark now that

$$(8) \quad \begin{aligned} \mathbb{E}_{\mathbf{H}}(\mathbf{1}_{\mathbf{x}}\mathbf{1}_{\mathbf{y}}) &= \mathbb{P}_{\mathbf{H}}(\mathbf{x}\mathbf{H}^{\top} = \mathbf{s} \text{ and } \mathbf{y}\mathbf{H}^{\top} = \mathbf{s}) \\ &= \mathbb{P}_{\mathbf{h}}(\varphi(\mathbf{h}) = (a, b))^{n-k} \end{aligned}$$

where in the last line we used that each rows of  $\mathbf{H}$  are independent and uniformly distributed. To conclude the proof it remains to plug Equation (7) in Equation (8).  $\square$

Lemma 5 enables us to deduce that

$$(9) \quad \begin{aligned} \sum_{\substack{\mathbf{x}, \mathbf{y} \in \mathcal{S}_t \\ \mathbf{x} \neq \mathbf{y}}} \mathbb{E}_{\mathbf{H}}(\mathbf{1}_{\mathbf{x}}\mathbf{1}_{\mathbf{y}}) - \mathbb{E}_{\mathbf{H}}(\mathbf{1}_{\mathbf{x}})\mathbb{E}_{\mathbf{H}}(\mathbf{1}_{\mathbf{y}}) &\leq \sum_{\mathbf{x} \in \mathcal{S}_t} \sum_{\substack{\mathbf{y} \in \mathcal{S}_t \setminus \mathbf{x}: \\ \text{colinear to } \mathbf{x}}} \frac{1}{q^{n-k}} - \frac{1}{q^{2(n-k)}} \\ &\leq \sum_{\mathbf{x} \in \mathcal{S}_t} \sum_{\substack{\mathbf{y} \in \mathcal{S}_t \setminus \mathbf{x}: \\ \text{colinear to } \mathbf{x}}} \frac{1}{q^{n-k}} \\ &\leq \frac{(q-2)\binom{n}{t}(q-1)^t}{q^{n-k}} \end{aligned}$$

It gives by plugging (9) in (6)

$$\begin{aligned} \mathbb{P}_{\mathbf{H}}\left(\left|N_t(\mathcal{C}, \mathbf{s}) - \frac{\binom{n}{t}(q-1)^t}{q^{n-k}}\right| \geq a\right) &\leq \frac{1}{a^2} \left( \frac{\binom{n}{t}(q-1)^t}{q^{n-k}} + \frac{(q-2)\binom{n}{t}(q-1)^t}{q^{n-k}} \right) \\ &= \frac{(q-1)\binom{n}{t}(q-1)^t}{a^2 q^{n-k}} \end{aligned}$$

which concludes the proof.  $\square$

The point of this proposition is that, for relative weights  $t/n \in (\tau^-, \tau^+)$  the term  $\binom{n}{t}(q-1)^t/q^{n-k}$ , is exponentially large. Therefore, by carefully setting  $a$  in this case we deduce  $N_t(\mathcal{C}, \mathbf{s})$  with a very good precision. For instance, if  $t/n \in (\tau^-, \tau^+)$ , let  $a = \left(\frac{\binom{n}{t}(q-1)^t}{q^{n-k}}\right)^{3/4}$  to obtain:

$$\mathbb{P}_{\mathbf{H}}\left(\left|N_t(\mathcal{C}, \mathbf{s}) - \frac{\binom{n}{t}(q-1)^t}{q^{n-k}}\right| \geq \left(\frac{\binom{n}{t}(q-1)^t}{q^{n-k}}\right)^{3/4}\right) \leq (q-1) \sqrt{\frac{q^{n-k}}{\binom{n}{t}(q-1)^t}} \in \text{negl}(n).$$

The number of errors of weight  $t$  that reach some syndrome is with an overwhelming probability equal to its expectation  $\binom{n}{t}(q-1)^t/q^{n-k}$  up to an additive factor  $(\binom{n}{t}(q-1)^t/q^{n-k})^{3/4}$  which is exponentially small with respect to  $\binom{n}{t}(q-1)^t/q^{n-k}$ .

#### 4. EXPECTED MINIMUM DISTANCE OF CODES

We are now interested in computing the expected minimum distance of a random code. As we will see it is given by the so-called *Gilbert-Varshamov* distance. This result is for cryptographic purposes very important. Suppose that one finds in a code a word of Hamming weight much smaller than it is expected. This would mean that the code is peculiar and maybe even worse, this codeword of small weight may reveal some secret information.

In view of the foregoing, it may be tempting to say that the expected minimum distance of a random  $[n, k]_q$ -code is given by the largest  $t$  such that  $\sum_{\ell \leq t} \frac{\binom{n}{\ell} (q-1)^\ell}{q^{n-k}} \leq 1$  and that is indeed what happens. It turns out that this value of  $t$  plays an important role in coding theory and is known as the Gilbert-Varshamov distance.

**Definition 1** (Gilbert-Varshamov Distance). *Let  $k \leq n$  and  $q$  be integers. The Gilbert-Varshamov distance  $t_{\text{GV}}(q, n, k)$  is defined as the largest integer such that*

$$\sum_{\ell=0}^{t_{\text{GV}}(q, n, k)} \binom{n}{\ell} (q-1)^\ell \leq q^{n-k}.$$

**Remark 5.** *The Gilbert-Varshamov distance gives the maximum  $r$  such that the volume of the ball of radius  $r$  is smaller than the inverse density of any  $[n, k]_q$ -code. Its analogue for lattices (in the context of lattice-based cryptography) is called the Gaussian heuristic.*

It can be verified (by using Lemma 1)

$$\frac{t_{\text{GV}}(q, n, Rn)}{n} = \tau^- + o(1)$$

where  $\tau^-$  is defined in Equation (4). It explains why  $\tau^-$  is commonly called the *relative Gilbert-Varshamov distance*. In the following proposition we show that the minimum distance of almost all  $[n, k]_q$ -codes is given by  $\tau^-$ . Interestingly, the proof that  $d_{\min}(\mathcal{C})/n > \tau^-$  happens with a negligible probability relies on Proposition 3 that used the second moment technique.

**Proposition 4.** *Let  $\varepsilon > 0$ . We have*

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

where  $\alpha \stackrel{\text{def}}{=} \min((1 - R) - h_q((1 + \varepsilon)\tau^-), h_q((1 - \varepsilon)\tau^-) - (1 - R)) > 0$ .

*Proof.* First notice that,

$$(10) \quad \mathbb{P}_{\mathbf{H}} \left( \frac{d_{\min}(\mathcal{C})}{n} \notin ((1 - \varepsilon)\tau^-, (1 + \varepsilon)\tau^-) \right) \leq \mathbb{P}_{\mathbf{H}} \left( \frac{d_{\min}(\mathcal{C})}{n} \leq (1 - \varepsilon)\tau^- \right) + \mathbb{P}_{\mathbf{H}} \left( \frac{d_{\min}(\mathcal{C})}{n} \geq (1 + \varepsilon)\tau^- \right).$$

Let us upper-bound independently both terms of the above inequation. First,

$$\begin{aligned} \mathbb{P}_{\mathbf{H}} \left( \frac{d_{\min}(\mathcal{C})}{n} \leq (1 - \varepsilon)\tau^- \right) &= \mathbb{P}_{\mathbf{H}} \left( \exists \mathbf{x} : \frac{|\mathbf{x}|}{n} \leq (1 - \varepsilon)\tau^- \text{ and } \mathbf{x} \in \mathcal{C} \right) \\ &\leq \sum_{\substack{\mathbf{x}: \\ \frac{|\mathbf{x}|}{n} \leq (1 - \varepsilon)\tau^-}} \mathbb{P}_{\mathbf{H}}(\mathbf{x} \in \mathcal{C}) \\ &= \sum_{\ell=0}^{(1 - \varepsilon)n\tau^-} \frac{\binom{n}{\ell} (q-1)^\ell}{q^{n-k}} \quad (\text{By Lemma 4.}) \\ &= \frac{q^{n(h_q((1 - \varepsilon)\tau^-) + o(1))}}{q^{n-k}} \end{aligned} \quad (11)$$

$$(12) \quad = q^{n(h_q((1 - \varepsilon)\tau^-) - (1 - R) + o(1))}$$

where we used in (11) Lemma 1 and the fact that  $x \mapsto h_q(x)$  is an increasing function for  $x \in [0, \tau^-] \subseteq \left[0, \frac{q-1}{q}\right]$ .

Let us now upper-bound the second term of (10). Let  $\delta \in (0, \varepsilon)$  and  $u \stackrel{\text{def}}{=} (1 + \delta)n\tau^-$ . Notice now that,

$$\begin{aligned}
 \mathbb{P}_{\mathbf{H}} \left( \frac{d_{\min}(\mathcal{C})}{n} \geq (1 + \varepsilon)\tau^- \right) &\leq \mathbb{P}_{\mathbf{H}} (N_u(\mathcal{C}, \mathbf{0}) = 0) \\
 &\leq \mathbb{P}_{\mathbf{H}} \left( \left| N_u(\mathcal{C}, \mathbf{0}) - \frac{\binom{n}{u}(q-1)^u}{q^{n-k}} \right| \geq \frac{\binom{n}{u}(q-1)^u}{q^{n-k}} \right) \\
 (13) \quad &\leq (q-1) \frac{q^{n-k}}{\binom{n}{u}(q-1)^u}
 \end{aligned}$$

where in the last line we used Proposition 3 by setting  $a = \frac{\binom{n}{u}(q-1)^u}{q^{n-k}}$ . But now by using Lemma 1 and that  $u = (1 + \delta)\tau^-$  we obtain:

$$\frac{q^{n-k}}{\binom{n}{u}(q-1)^u} = q^{n(1-R-h_q((1+\delta)\tau^-)+o(1))}.$$

By plugging this in Equation (13) we obtain for any  $\delta \in (0, \varepsilon)$ :

$$\mathbb{P}_{\mathbf{H}} \left( \frac{d_{\min}(\mathcal{C})}{n} \geq (1 + \varepsilon)\tau^- \right) \leq q^{n(1-R-h_q((1+\delta)\tau^-)+o(1))}.$$

Therefore, by letting  $\delta \rightarrow \varepsilon$  we get:

$$(14) \quad \mathbb{P}_{\mathbf{H}} \left( \frac{d_{\min}(\mathcal{C})}{n} \geq (1 + \varepsilon)\tau^- \right) \leq q^{n(1-R-h_q((1+\varepsilon)\tau^-)+o(1))}.$$

To conclude the proof it remains to put together Equations (12) and (14) in Equation (10).  $\square$

**Remark 6.** In coding theory, the relative Gilbert-Varshamov distance  $\tau^-$  is known as a “lower-bound”, there exists a family of  $[n, Rn]_q$ -codes with a relative minimum distance  $\geq \tau^-$ . What we have actually proven is that almost all families of  $[n, Rn]_q$ -codes have asymptotically a relative minimum distance  $= \tau^-$ . Let us stress that it does not show that the relative Gilbert-Varshamov distance  $\tau^-$  is an “upper-bound”, all families of  $[n, Rn]_q$ -codes are such that their asymptotically relative minimum distance is  $\leq \tau^-$ . This is a widely open conjecture in the case of  $\mathbb{F}_2$  but which is not true for  $\mathbb{F}_q$  as long as  $q$  is a square  $\geq 49$ .

**About the optimality of random codes.** We have seen in lecture notes 1 that balls centred at codewords  $\mathbf{c} \in \mathcal{C}$  and whose radius is  $< \frac{d_{\min}(\mathcal{C})}{2}$  never overlap. This condition over the radius ensures that when an error  $\mathbf{e}$  of Hamming weight smaller than  $< \frac{d_{\min}(\mathcal{C})}{2}$  occurs, we are sure that computing the closest codeword from  $\mathbf{c} + \mathbf{e}$  will outcome  $\mathbf{c}$  (we say that the *maximum likelihood decoding* succeeds). However, for random codes this property can be made even stronger. In that case we can show that  $\mathbf{c}$  is indeed the closest codeword from  $\mathbf{c} + \mathbf{e}$  (with overwhelming probability) if  $\mathbf{e} \approx d_{\min}(\mathcal{C})$  where  $d_{\min}(\mathcal{C}) \approx \tau^-n$  as we show now in the following proposition (by using Markov’s inequality).

**Proposition 5.** Let  $t \stackrel{\text{def}}{=} n(1 - \varepsilon)\tau^-$  for some  $\varepsilon > 0$  and let  $\eta > 0$ . We have

$$\mathbb{P}_{\mathbf{H}} \left( \frac{\#\{\mathbf{c}, \mathbf{c}' \in \mathcal{C}, \mathbf{e}, \mathbf{e}' \in \mathcal{S}_t : \mathbf{c} + \mathbf{e} = \mathbf{c}' + \mathbf{e}'\}}{q^k \binom{n}{t} (q-1)^t} \geq 1 + \eta \right) \leq \frac{1}{\eta} \frac{\binom{n}{t} (q-1)^t}{q^{n-k}} = \frac{1}{\eta} q^{-\alpha n(1+o(1))}$$

where  $\alpha \stackrel{\text{def}}{=} h_q((1 - \varepsilon)\tau^-) - 1 + R > 0$ .

Notice that,

$$\frac{\#\{\mathbf{c}, \mathbf{c}' \in \mathcal{C}, \mathbf{e}, \mathbf{e}' \in \mathcal{S}_t : \mathbf{c} + \mathbf{e} = \mathbf{c}' + \mathbf{e}'\}}{q^k \binom{n}{t} (q-1)^t} = 1$$

means that there are no collisions between noisy codewords with errors of Hamming weight  $t$  (and that balls of radius  $t$  and centered at codewords do not overlap). Therefore, the above proposition shows that

for random codes, the maximum likelihood decoding will succeed for almost all noisy codewords, up to the Gilbert-Varshamov distance if  $\eta$  is chosen sufficiently small. However, once again, we could be more accurate by using the second moment technique with Bienaymé-Tchebychev's inequality.

*Proof.* Let,

$$Z \stackrel{\text{def}}{=} \# \{ \mathbf{c}, \mathbf{c}' \in \mathcal{C}, \mathbf{e}, \mathbf{e}' \in \mathcal{S}_t : \mathbf{c} + \mathbf{e} = \mathbf{c}' + \mathbf{e}' \}.$$

We have the following computation

$$\begin{aligned} Z &= \sum_{\mathbf{c} \in \mathcal{C}, \mathbf{e} \in \mathcal{S}_t} 1 + \sum_{\substack{\mathbf{c}, \mathbf{c}' \in \mathcal{C}, \mathbf{e}, \mathbf{e}' \in \mathcal{S}_t \\ (\mathbf{c}, \mathbf{e}) \neq (\mathbf{c}', \mathbf{e}') \\ \mathbf{c} + \mathbf{e} = \mathbf{c}' + \mathbf{e}'}} 1 \\ &= q^k \binom{n}{t} (q-1)^t + q^k \sum_{\substack{\mathbf{e}, \mathbf{e}' \in \mathcal{S}_t \\ \mathbf{e} \neq \mathbf{e}' : (\mathbf{e} - \mathbf{e}') \mathbf{H}^\top = \mathbf{0}}} 1 \\ &= q^k \binom{n}{t} (q-1)^t \left( 1 + \frac{1}{\binom{n}{t} (q-1)^t} \sum_{\substack{\mathbf{e}, \mathbf{e}' \in \mathcal{S}_t \\ \mathbf{e} \neq \mathbf{e}' : (\mathbf{e} - \mathbf{e}') \mathbf{H}^\top = \mathbf{0}}} 1 \right). \end{aligned}$$

Let,

$$X \stackrel{\text{def}}{=} \frac{1}{\binom{n}{t} (q-1)^t} \sum_{\substack{\mathbf{e}, \mathbf{e}' \in \mathcal{S}_t \\ \mathbf{e} \neq \mathbf{e}' : (\mathbf{e} - \mathbf{e}') \mathbf{H}^\top = \mathbf{0}}} 1$$

By using Lemma 4,

$$\begin{aligned} \mathbb{E}_{\mathbf{H}}(X) &= \frac{1}{\binom{n}{t} (q-1)^t} \sum_{\substack{\mathbf{e}, \mathbf{e}' \in \mathcal{S}_t \\ \mathbf{e} \neq \mathbf{e}'}} \frac{1}{q^{n-k}} \\ &\leq \frac{1}{\binom{n}{t} (q-1)^t} \frac{\left( \binom{n}{t} (q-1)^t \right)^2}{q^{n-k}} \\ &= \frac{\binom{n}{t} (q-1)^t}{q^{n-k}}. \end{aligned}$$

To conclude the proof it is enough to apply Markov's inequality to upper-bound  $\mathbb{P}_{\mathbf{H}}(X > \eta)$ .  $\square$

## 5. UNIFORM DISTRIBUTION OF SYNDROMES

In lecture notes 1 we have seen that the decoding problem is defined (for cryptographic purposes) with some distribution in input, namely  $(\mathbf{H}, \mathbf{xH}^\top)$  where  $\mathbf{H} \in \mathbb{F}_q^{(n-k) \times n}$  and  $\mathbf{x} \in \mathcal{S}_t$  are uniformly distributed. Our aim in what follows is to show that when  $t/n \in (\tau^-, \tau^+)$  we could replace  $\mathbf{xH}^\top$  by  $\mathbf{s} \in \mathbb{F}_q^{n-k}$  being uniformly distributed, without changing our problem. More precisely, we are going to show that distributions  $(\mathbf{H}, \mathbf{xH}^\top)$  and  $(\mathbf{H}, \mathbf{s})$  are *statistically close* under the condition  $t/n \in (\tau^-, \tau^+)$ , showing that for any algorithm solving  $\text{DP}(n, q, R, \tau)$  we could choose its inputs as  $(\mathbf{H}, \mathbf{s})$  without changing at much its success probability. This result may seem useless but in some applications where  $\tau \in (\tau^-, \tau^+)$  it is much more comfortable to consider directly  $(\mathbf{H}, \mathbf{s})$  rather than  $(\mathbf{H}, \mathbf{xH}^\top)$ .

**Proposition 6.** Let  $k \stackrel{\text{def}}{=} \lfloor Rn \rfloor$ ,  $t \stackrel{\text{def}}{=} \lfloor \tau n \rfloor$  and  $\mathbf{H} \in \mathbb{F}_q^{(n-k) \times n}$ ,  $\mathbf{s} \in \mathbb{F}_q^{n-k}$ ,  $\mathbf{e} \in \mathcal{S}_t$  being uniformly distributed. We have

$$(15) \quad \mathbb{E}_{\mathbf{H}} (\Delta (\mathbf{eH}^\top, \mathbf{s})) \leq \frac{1}{2} \sqrt{\frac{q^{n-k} - 1}{\binom{n}{t}(q-1)^t}}.$$

In particular, when  $\tau \in (\tau^-, \tau^+)$

$$\mathbb{E}_{\mathbf{H}} (\Delta (\mathbf{eH}^\top, \mathbf{s})) \leq q^{-\alpha n(1+o(1))} \quad \text{where} \quad \alpha \stackrel{\text{def}}{=} \frac{1}{2} (h_q(\tau) - 1 + R) > 0.$$

**Exercise 5.** Let  $\mathbf{H} \in \mathbb{F}_q^{(n-k) \times n}$  being uniformly distributed,  $\mathbf{s} \in \mathbb{F}_q^{n-k}$ ,  $\mathbf{e} \in \mathcal{S}_t$  be some random variables. Show that

$$\mathbb{E}_{\mathbf{H}} (\Delta (\mathbf{eH}^\top, \mathbf{s})) = \frac{1}{q^{(n-k) \times n}} \sum_{\mathbf{H}_0 \in \mathbb{F}_q^{(n-k) \times n}} \Delta (\mathbf{eH}_0^\top, \mathbf{s})$$

The proof of Proposition 6 relies on the following lemma which is a rewriting in our context of the result known as the *left over hash lemma*. Roughly speaking, it shows that if the outputs of a function collide with a probability  $\varepsilon$ -close to the case where they would be randomly distributed, then this function is  $\sqrt{\varepsilon}$ -close to a random function.

**Lemma 6.** Let  $\mathcal{H} = (h_i)_{i \in I}$  be a finite family of applications from  $E$  in  $F$ . Let  $\varepsilon$  be the “collision bias”

$$\mathbb{P}_{h,e,e'}(h(e) = h(e')) = \frac{1}{\#F}(1 + \varepsilon)$$

where  $h$  is uniformly drawn in  $\mathcal{H}$ ,  $e$  and  $e'$  be distributed according to some random variable  $X$  taking its values  $E$ . Let  $\mathcal{U}$  be the uniform distribution over  $F$  and  $\mathcal{D}(h)$  be the distribution  $h(e)$  when  $e$  is distributed according to  $X$ . We have,

$$\mathbb{E}_h (\Delta(\mathcal{D}(h), \mathcal{U})) \leq \frac{1}{2} \sqrt{\varepsilon}.$$

*Proof.* By definition of the statistical distance we have

$$(16) \quad \begin{aligned} \mathbb{E}_h (\Delta(\mathcal{D}(h), \mathcal{U})) &= \sum_{h \in \mathcal{H}} \frac{1}{\#\mathcal{H}} \Delta(\mathcal{D}(h), \mathcal{U}) \\ &= \frac{1}{2} \sum_{h \in \mathcal{H}} \frac{1}{\#\mathcal{H}} \sum_{f \in F} \left| \mathbb{P}_e(h(e) = f) - \frac{1}{\#F} \right| \\ &= \frac{1}{2} \sum_{(h,f) \in \mathcal{H} \times F} \left| \mathbb{P}_{h_0,e}(h_0 = h, h_0(e) = f) - \frac{1}{\#\mathcal{H} \#F} \right| \\ (17) \quad &= \frac{1}{2} \sum_{(h,f) \in \mathcal{H} \times F} \left| q_{h,f} - \frac{1}{\#\mathcal{H} \#F} \right|. \end{aligned}$$

where  $q_{h,f} \stackrel{\text{def}}{=} \mathbb{P}_{h_0,e}((h_0, h_0(e)) = (h, f))$  with  $h_0$  being uniformly chosen at random in  $\mathcal{H}$  and  $e$  be distributed according to  $X$ . Using the Cauchy-Schwarz inequality, we obtain

$$(18) \quad \sum_{(h,f) \in \mathcal{H} \times F} \left| q_{h,f} - \frac{1}{\#\mathcal{H} \#F} \right| \leq \sqrt{\sum_{(h,f) \in \mathcal{H} \times F} \left( q_{h,f} - \frac{1}{\#\mathcal{H} \#F} \right)^2} \sqrt{\#\mathcal{H} \#F}.$$

Let us observe now that

$$\begin{aligned}
 \sum_{(h,f) \in \mathcal{H} \times F} \left( q_{h,f} - \frac{1}{\#\mathcal{H} \#F} \right)^2 &= \sum_{h,f} \left( q_{h,f}^2 - 2 \frac{q_{h,f}}{\#\mathcal{H} \#F} + \frac{1}{\#\mathcal{H}^2 \#F^2} \right) \\
 &= \sum_{h,f} q_{h,f}^2 - 2 \frac{\sum_{h,f} q_{h,f}}{\#\mathcal{H} \#F} + \frac{1}{\#\mathcal{H} \#F} \\
 &= \sum_{h,f} q_{h,f}^2 - \frac{1}{\#\mathcal{H} \#F}.
 \end{aligned}
 \tag{19}$$

Consider for  $i \in \{0, 1\}$  independent random variables  $h_i$  and  $e_i$  that are drawn uniformly at random in  $\mathcal{H}$  and according to  $X$  respectively. We continue this computation by noticing now that

$$\begin{aligned}
 \sum_{h,f} q_{h,f}^2 &= \sum_{h,f} \mathbb{P}_{h_0, e_0}(h_0 = h, h_0(e_0) = f) \mathbb{P}_{h_1, e_1}(h_1 = h, h_1(e_1) = f) \\
 &= \mathbb{P}_{h_0, h_1, e_0, e_1}(h_0 = h_1, h_0(e_0) = h_1(e_1)) \\
 &= \frac{\mathbb{P}_{h_0, e_0, e_1}(h_0(e_0) = h_0(e_1))}{\#\mathcal{H}} \\
 &= \frac{1 + \varepsilon}{\#\mathcal{H} \#F}.
 \end{aligned}
 \tag{20}$$

By substituting for  $\sum_{h,f} q_{h,f}^2$  the expression obtained in (20) into (19) and then back into (18) we finally obtain

$$\begin{aligned}
 \sum_{(h,f) \in \mathcal{H} \times F} \left| q_{h,f} - \frac{1}{\#\mathcal{H} \#F} \right| &\leq \sqrt{\frac{1 + \varepsilon}{\#\mathcal{H} \#F} - \frac{1}{\#\mathcal{H} \#F}} \sqrt{\#\mathcal{H} \#F} \\
 &= \sqrt{\frac{\varepsilon}{\#\mathcal{H} \#F}} \sqrt{\#\mathcal{H} \#F} \\
 &= \sqrt{\varepsilon}.
 \end{aligned}$$

□

We are now ready to prove Proposition 6.

*Proof of Proposition 6.* Let us compute the “collision bias” in our case. We have:

$$\begin{aligned}
 \mathbb{P}_{\mathbf{H}, \mathbf{e}, \mathbf{e}'}(\mathbf{e}\mathbf{H}^\top = \mathbf{e}'\mathbf{H}^\top) &= \mathbb{P}_{\mathbf{H}, \mathbf{e}, \mathbf{e}'}((\mathbf{e} - \mathbf{e}')\mathbf{H}^\top = \mathbf{0}) \\
 &= \mathbb{P}_{\mathbf{H}, \mathbf{e}, \mathbf{e}'}((\mathbf{e} - \mathbf{e}')\mathbf{H}^\top = \mathbf{0} \mid \mathbf{e} \neq \mathbf{e}') \mathbb{P}(\mathbf{e} \neq \mathbf{e}') + \mathbb{P}_{\mathbf{e}, \mathbf{e}'}(\mathbf{e} = \mathbf{e}') \\
 &= \frac{1}{q^{n-k}} \left( 1 - \frac{1}{\binom{n}{t}(q-1)^t} \right) + \frac{1}{\binom{n}{t}(q-1)^t} \quad (\text{by Lemma 4}) \\
 &= \frac{1}{q^{n-k}} (1 + \varepsilon) \quad \text{where} \quad \varepsilon \stackrel{\text{def}}{=} \frac{q^{n-k} - 1}{\binom{n}{t}(q-1)^t}
 \end{aligned}$$

which shows (15) by using Lemma 6. The second part of the proposition easily follows from the definition of  $\tau^-$  and  $\tau^+$ . □

**Exercise 6.** Show that if one replaces in the binary case ( $q = 2$ )  $\mathbf{e}$  in Proposition 6 by  $\mathbf{e}^{\text{Ber}}$ , where the  $e_i^{\text{Ber}}$ ’s are distributed according to a Bernoulli distribution of parameter  $\tau$ , we would obtain

$$\mathbb{E}_{\mathbf{H}}(\Delta(\mathbf{e}^{\text{Ber}}\mathbf{H}^\top, \mathbf{s})) \leq \frac{1}{2} \sqrt{2^{-k} (1 + (1 - 2\tau)^2)^n}.$$

What can you deduce when comparing both results with  $\mathbf{e}$  or  $\mathbf{e}^{\text{Ber}}$ ? What is (according to Proposition 6) the “best” choice of error  $\mathbf{x}$  to ensure that  $\mathbf{x}\mathbf{H}^\top$  is uniformly distributed?

**Exercise 7.** Let  $\mathcal{C}$  be a fixed  $[n, k]_q$ -code of parity-check matrix  $\mathbf{H}$  and  $\mathbf{y}, \mathbf{s}, \mathbf{e} \in \mathbb{F}_q^n \times \mathbb{F}_q^{n-k} \times \mathcal{S}_t$  be uniformly distributed. Our aim in this exercise is to show that  $\Delta(\mathbf{c} + \mathbf{e}, \mathbf{y}) = \Delta(\mathbf{e}\mathbf{H}^\top, \mathbf{s})$ .

1. Given  $\mathbf{s} \in \mathbb{F}_q^{n-k}$ , let  $\mathbf{y}(\mathbf{s}) \in \mathbb{F}_q^n$  be such that  $\mathbf{y}(\mathbf{s})\mathbf{H}^\top = \mathbf{s}$ . Show that

$$\sum_{\mathbf{y} \in \mathbb{F}_q^n} \left| \mathbb{P}_{\mathbf{e}, \mathbf{c}}(\mathbf{c} + \mathbf{e} = \mathbf{y}) - \frac{1}{q^n} \right| = \sum_{\mathbf{s} \in \mathbb{F}_q^{n-k}} \sum_{\mathbf{c}' \in \mathcal{C}} \left| \mathbb{P}_{\mathbf{e}, \mathbf{c}}(\mathbf{c} + \mathbf{e} = \mathbf{y}(\mathbf{s}) + \mathbf{c}') - \frac{1}{q^n} \right|.$$

2. Deduce that  $\Delta(\mathbf{c} + \mathbf{e}, \mathbf{y}) = \Delta(\mathbf{e}\mathbf{H}^\top, \mathbf{s})$ .

Up to now, we have proven in Proposition 6 that syndromes  $\mathbf{e}\mathbf{H}^\top$  are statistically close to the uniform distribution when  $\mathbf{e}$  is picked uniformly at random in  $\mathcal{S}_t$  with  $t$  larger than the Gilbert-Varshamov distance (more precisely, when  $t/n \in (\tau^-, \tau^+)$ ) but in average over  $\mathbf{H}$ . One may ask if the result still holds for a fixed matrix  $\mathbf{H}_0$ ? Actually we can prove, without too much effort, that the above result is true for almost all matrices but with a loss given by a square root as shown by the following proposition.

**Proposition 7.** Let  $k \stackrel{\text{def}}{=} \lfloor Rn \rfloor$ ,  $t \stackrel{\text{def}}{=} \lfloor \tau n \rfloor$ ,  $\mathbf{H} \in \mathbb{F}_q^{(n-k) \times n}$  being uniformly distributed and  $\mathbf{s} \in \mathbb{F}_q^{n-k}$ ,  $\mathbf{e} \in \mathbb{F}_q^n$  be some random variables. Suppose that

$$\mathbb{E}_{\mathbf{H}} (\Delta(\mathbf{e}\mathbf{H}^\top, \mathbf{s})) \leq \varepsilon$$

Then, we have

$$\frac{\#\left\{ \mathbf{H}_0 \in \mathbb{F}_q^{(n-k) \times n} : \Delta(\mathbf{e}\mathbf{H}_0^\top, \mathbf{s}) \geq \sqrt{\varepsilon} \right\}}{q^{(n-k) \times n}} \leq \sqrt{\varepsilon}$$

*Proof.* First,

$$(21) \quad \mathbb{E}_{\mathbf{H}} (\Delta(\mathbf{e}\mathbf{H}^\top, \mathbf{s})) = \frac{1}{q^{(n-k) \times n}} \sum_{\mathbf{H}_0 \in \mathbb{F}_q^{(n-k) \times n}} \Delta(\mathbf{e}\mathbf{H}_0^\top, \mathbf{s}).$$

The idea of the proof is to split in two parts this sum according to the terms that are larger and smaller than  $\sqrt{\varepsilon}$

$$\begin{aligned} \sum_{\mathbf{H}_0 \in \mathbb{F}_q^{(n-k) \times n}} \Delta(\mathbf{e}\mathbf{H}_0^\top, \mathbf{s}) &= \sum_{\substack{\mathbf{H}_0 : \\ \Delta(\mathbf{e}\mathbf{H}_0^\top, \mathbf{s}) < \sqrt{\varepsilon}}} \Delta(\mathbf{e}\mathbf{H}_0^\top, \mathbf{s}) + \sum_{\substack{\mathbf{H}_0 : \\ \Delta(\mathbf{e}\mathbf{H}_0^\top, \mathbf{s}) \geq \sqrt{\varepsilon}}} \Delta(\mathbf{e}\mathbf{H}_0^\top, \mathbf{s}) \\ &\geq \sum_{\substack{\mathbf{H}_0 : \\ \Delta(\mathbf{e}\mathbf{H}_0^\top, \mathbf{s}) < \sqrt{\varepsilon}}} \Delta(\mathbf{e}\mathbf{H}_0^\top, \mathbf{s}) + \sum_{\substack{\mathbf{H}_0 : \\ \Delta(\mathbf{e}\mathbf{H}_0^\top, \mathbf{s}) \geq \sqrt{\varepsilon}}} \sqrt{\varepsilon} \\ &\geq \sum_{\substack{\mathbf{H}_0 : \\ \Delta(\mathbf{e}\mathbf{H}_0^\top, \mathbf{s}) \geq \sqrt{\varepsilon}}} \sqrt{\varepsilon} \\ (22) \quad &= \sqrt{\varepsilon} \#\left\{ \mathbf{H}_0 \in \mathbb{F}_q^{(n-k) \times n} : \Delta(\mathbf{e}\mathbf{H}_0^\top, \mathbf{s}) \geq \sqrt{\varepsilon} \right\} \end{aligned}$$

By plugging Equation (22) in (21), we obtain

$$\mathbb{E}_{\mathbf{H}} (\Delta(\mathbf{e}\mathbf{H}^\top, \mathbf{s})) \geq \sqrt{\varepsilon} \frac{\#\left\{ \mathbf{H}_0 \in \mathbb{F}_q^{(n-k) \times n} : \Delta(\mathbf{e}\mathbf{H}_0^\top, \mathbf{s}) \geq \sqrt{\varepsilon} \right\}}{q^{(n-k) \times n}}$$

But by assumption we have  $\mathbb{E}_{\mathbf{H}} (\Delta(\mathbf{e}\mathbf{H}^\top, \mathbf{s})) \leq \varepsilon$ . It concludes the proof.  $\square$