

6.7720/18.619/15.070 Lecture 5

Shannon's Noisy Coding Theorem

Kuikui Liu

February 18, 2025

Acknowledgements & Disclaimers *In the process of writing these notes, we consulted materials by Venkatesan Guruswami and Michel X. Goemans. Please be advised that these notes have not been subjected to the usual scrutiny reserved for formal publications. If you do spot an error, please contact the instructor.*

1 Error-Correction and Shannon's Theorem

In this lecture, we use the probabilistic method combined with the Weak Law of Large Numbers to prove a remarkable characterization of the fundamental limits of error-correcting mechanisms for communication in the presence of noise. We will work with the following bare-bones definition of an error-correcting code.

Definition 1 (Error-Correcting Code). *For positive integers $k, n \in \mathbb{N}$ satisfying $n \geq k$, an (error-correcting) $[n, k]$ -code is a pair of functions $\text{Enc} : \{0, 1\}^k \rightarrow \{0, 1\}^n$ and $\text{Dec} : \{0, 1\}^n \rightarrow \{0, 1\}^k$. We will often refer to the former as the encoder, and the latter as the decoder. The rate of the code is the ratio $r = \frac{k}{n}$. We often refer to the elements of $\{0, 1\}^k$ as messages, and the elements of $\{0, 1\}^n$ as codewords.*

One should imagine that a sender would like to transmit a message $\mathbf{m} \in \{0, 1\}^k$ to a receiver, but only has access to a noisy communication channel. Hence, the sender will first encode the message, and then send $\text{Enc}(\mathbf{m})$ through the channel. The hope is that even though the codeword $\text{Enc}(\mathbf{m})$ can suffer errors during transmission, the receiver will be able to recover \mathbf{m} by applying the decoder Dec to the (possibly corrupted) codeword they received. In practice, we typically want our error-correcting codes to have additional properties, e.g. *linearity, locality, efficient computability* of Enc, Dec , etc. We will not be concerned with these in this lecture.

To make the problem well-defined, we also need to fix a model for how errors will be introduced. In this lecture, we will only consider a very special type of error, namely every bit of the codeword can be flipped independently with some probability $0 \leq p \leq 1/2$. This *channel* should be familiar; we saw it in the context of the broadcast process on trees.

Definition 2 (Binary Symmetric Channel). *For $0 \leq p \leq 1/2$, the binary symmetric channel with parameter p , denoted BSC_p , is described by the following randomized operation: For any bit-string $\mathbf{x} \in \{0, 1\}^*$ of any length n , $\mathbf{y} = \text{BSC}_p(\mathbf{x})$ is a random length- n bit-string such that independently for each $i = 1, \dots, n$, $\mathbf{y}_i = \mathbf{x}_i$ with probability $1 - p$ and $\mathbf{y}_i \neq \mathbf{x}_i$ with probability p .*

One can, of course, consider more exotic error models (e.g. random erasures, random swaps, correlated bit flips, etc.), as well as *worst-case* error models.

For an error-correcting code (Enc, Dec) , we will refer to

$$\Pr[\text{Dec}(\text{BSC}_p(\text{Enc}(\mathbf{m}))) \neq \mathbf{m}] \quad \text{and} \quad \Pr[\text{Dec}(\text{BSC}_p(\text{Enc}(\mathbf{m}))) = \mathbf{m}]$$

as the *failure/error* and *success* probabilities of the message \mathbf{m} , respectively. The former (resp. latter) is the probability that if we send the encoded message $\text{Enc}(\mathbf{m})$ through the noisy BSC_p channel, then the receiver will be unable (resp. able) to recover \mathbf{m} using the decoding function $\text{Dec}(\cdot)$.

Naturally, there is a tradeoff between the rate of an error-correcting code, which quantifies its efficiency, and the (worst-case) failure probability of any given message. For instance, if you allow

arbitrarily small rates (i.e. arbitrarily inefficient codes), then you can drive down the failure probabilities uniformly by using something like the *repetition code*, which simply copies each message bit some large number of times.

Remarkably, there turns out to be a specific threshold for the rate, above which every error-correcting code must fail for some message with probability close to 1, and below which there exists an error-correcting code with failure probabilities uniformly close to 0. This phase transition phenomenon is the content of Shannon's seminal Noisy Coding Theorem. To state it, define the *binary entropy function* by

$$H(p) \stackrel{\text{def}}{=} -p \log_2 p - (1-p) \log_2 (1-p), \quad \forall p \in [0, 1].$$

Theorem 1.1 (Shannon's Noisy Coding Theorem; [???]). *Fix any constant $0 \leq p < 1/2$.*

- *Suppose $r > 1 - H(p)$. Then for every $n \in \mathbb{N}$ and every $[n, k]$ -code (Enc, Dec) with $k = rn$,*

$$\max_{\mathbf{m} \in \{0,1\}^k} \Pr [\text{Dec}(\text{BSC}_p(\text{Enc}(\mathbf{m}))) \neq \mathbf{m}] \geq 1 - o(1).$$

- *Suppose $r < 1 - H(p)$. Then for every $n \in \mathbb{N}$, there exists an $[n, k]$ -code (Enc, Dec) with $k = rn$ such that*

$$\max_{\mathbf{m} \in \{0,1\}^k} \Pr [\text{Dec}(\text{BSC}_p(\text{Enc}(\mathbf{m}))) \neq \mathbf{m}] \leq o(1).$$

Remark 1. We note that Shannon's Noisy Coding Theorem actually holds for much more general channels beyond the BSC_p , although one must replace $1 - H(p)$ by a more general information-theoretic quantity known as the *channel capacity*.

To prove [Theorem 1.1](#), we will need an extremely weak concentration bound for sums of independent random variables, which can be obtained directly from the second moment. By leveraging more powerful concentration estimates like *Chernoff bounds* (which we introduce in the next lecture), one can obtain strong quantitative control on the $o(1)$ terms in [Theorem 1.1](#).

Theorem 1.2 (Weak Law of Large Numbers). *Let $\{X_i\}_{i=1}^\infty$ be a sequence of independent and identically distributed (i.i.d.) random variables with mean μ and variance $\sigma^2 < \infty$. Then*

$$\lim_{n \rightarrow \infty} \Pr \left[\left| \frac{1}{n} \sum_{i=1}^n X_i - \mu \right| < \epsilon \right] = 1, \quad \forall \epsilon > 0.$$

In other words, we have convergence in probability $\frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{P} \mu$ as $n \rightarrow \infty$.

Remark 2. As the proof shows, one can significantly relax the joint independence assumption to merely *pairwise independence*, i.e. for all $i \neq j$ and all x_i, x_j , $\Pr[X_i = x_i, X_j = x_j] = \Pr[X_i = x_i] \cdot \Pr[X_j = x_j]$. A direct calculation reveals this latter condition is sufficient to guarantee that the variance of the sum is equal to the sum of the variances.

Proof. By independence, we have $\text{Var}(\frac{1}{n} \sum_{i=1}^n X_i) = \frac{\sigma^2}{n}$. Hence, by Chebyshev's Inequality,

$$\Pr \left[\left| \frac{1}{n} \sum_{i=1}^n X_i - \mu \right| \geq \epsilon \right] \leq \frac{\text{Var}(\frac{1}{n} \sum_{i=1}^n X_i)}{\epsilon^2} = \frac{\sigma^2}{\epsilon^2} \cdot \frac{1}{n}.$$

□

Finally, to quantify the amount of errors introduced to a codeword, we define a metric on bit-strings known as *Hamming distance*:

$$d_H(\mathbf{x}, \mathbf{y}) \stackrel{\text{def}}{=} \|\mathbf{x} - \mathbf{y}\|_1 = \#\{i \in [n] : x_i \neq y_i\}, \quad \forall \mathbf{x}, \mathbf{y} \in \{0, 1\}^n.$$

We have the following bound on the volume of Hamming balls, similar to what we used when we looked at the connectivity threshold of Erdős-Rényi random graphs.

Lemma 1.3 (Volume of a Hamming Ball). *For $0 \leq p \leq 1/2$,*

$$\#\{\mathbf{x} \in \{0,1\}^n : d_H(\mathbf{0}, \mathbf{x}) \leq pn\} = \sum_{k=0}^{pn} \binom{n}{k} \leq 2^{H(p) \cdot n}.$$

Proof. The first identity is immediate. For the second identity, we have

$$1 = (p + (1-p))^n \geq \sum_{k=0}^{pn} \binom{n}{k} p^k (1-p)^{n-k} \geq p^{pn} (1-p)^{(1-p)n} \sum_{k=0}^{pn} \binom{n}{k} = 2^{-H(p) \cdot n} \sum_{k=0}^{pn} \binom{n}{k},$$

where in the second inequality, we used $0 \leq p \leq 1/2$ and the fact that the function $k \mapsto p^k (1-p)^{n-k}$ is monotone decreasing for $0 \leq k \leq pn$. Rearranging then yields the upper bound. \square

2 Proof of Theorem 1.1: The Case $r > 1 - H(p)$

We will go via the contrapositive, and show that if there exists an error-correcting code for which the success probability is uniformly lower bounded by some positive constant $\epsilon > 0$ for all possible messages $\mathbf{m} \in \{0,1\}^k$, then it must be that $r \leq 1 - H(p) + o(1)$. To achieve this, our goal will be to argue that any message \mathbf{m} with constant success probability must have at least $2^{(H(p)-o(1)) \cdot n}$ codewords $\mathbf{x} \in \{0,1\}^n$ which decode to \mathbf{m} , i.e. $\text{Dec}(\mathbf{x}) = \mathbf{m}$. Since the preimages of the decoder Dec partition $\{0,1\}^n$, this will imply a bound on the number of possible such messages \mathbf{m} . This is a classic example of a *volume argument*.

Let us now formalize our strategy.

Proposition 2.1. *For any constant $\epsilon > 0$ independent of n , if $\mathbf{m} \in \{0,1\}^k$ is a message with success probability at least ϵ , then*

$$\#\{\mathbf{x} \in \{0,1\}^n : \text{Dec}(\mathbf{x}) = \mathbf{m}\} \geq 2^{(H(p)-o(1)) \cdot n}.$$

Let us first quickly see how this can be used to show $r \leq 1 - H(p) + o(1)$ and complete the proof of the first case of Theorem 1.1. For $r \in [0,1]$, suppose there exists $\epsilon > 0$ such that for infinitely many $n \in \mathbb{N}$, there exist $\text{Enc} : \{0,1\}^k \rightarrow \{0,1\}^n$ and $\text{Dec} : \{0,1\}^n \rightarrow \{0,1\}^k$ with $k = rn$ such that

$$\max_{\mathbf{m} \in \{0,1\}^k} \Pr[\text{Dec}(\text{BSC}_p(\text{Enc}(\mathbf{m}))) \neq \mathbf{m}] \leq 1 - \epsilon$$

Then the success probability of every $\mathbf{m} \in \{0,1\}^k$ is uniformly lower bounded by ϵ . Hence,

$$\begin{aligned} 2^n &\geq \sum_{\mathbf{m} \in \{0,1\}^k} \#\{\mathbf{x} \in \{0,1\}^n : \text{Dec}(\mathbf{x}) = \mathbf{m}\} && \text{(Preimages of Dec partition } \{0,1\}^n) \\ &\geq 2^k \cdot 2^{(H(p)-o(1)) \cdot n}. && \text{(Invoking Proposition 2.1)} \end{aligned}$$

Rearranging then yields the inequality $r \leq 1 - H(p) + o(1)$. Since this holds for infinitely many n , we must have $r \leq 1 - H(p)$ as desired.

Proof of Proposition 2.1. Our assumption says that applying BSC_p to $\text{Enc}(\mathbf{m})$ yields a random bit-string which happens to fall in the preimage $\{\mathbf{x} \in \{0,1\}^n : \text{Dec}(\mathbf{x}) = \mathbf{m}\}$ with constant probability:

$$\Pr[\text{BSC}_p(\text{Enc}(\mathbf{m})) \in \text{Dec}^{-1}(\mathbf{m})] \geq \epsilon.$$

Our goal is to use this lower bound to deduce a lower bound on the cardinality of the preimage $\text{Dec}^{-1}(\mathbf{m})$. Observe that if $A \subseteq \{0,1\}^n$ is any subset of codewords, then

$$|\text{Dec}^{-1}(\mathbf{m})| \geq |\text{Dec}^{-1}(\mathbf{m}) \cap A| = \Pr_{\mathbf{x} \sim \text{Unif}(A)}[\mathbf{x} \in \text{Dec}^{-1}(\mathbf{m})] \cdot |A|.$$

We will employ an inequality of this flavor. Somewhat inconveniently, the probability in our assumption is with respect to a nonuniform distribution over $\{0,1\}^n$. However, we can get around this by using concentration. For convenience, let $\mathbf{y} = \text{Enc}(\mathbf{m})$. For $\delta > 0$, let A_δ denote the *annulus* around \mathbf{y}

$$A_\delta \stackrel{\text{def}}{=} \{\mathbf{x} \in \{0,1\}^n : d_H(\mathbf{x}, \mathbf{y}) \in [p - \delta, p + \delta]\}.$$

The idea is to choose $\delta(n) \rightarrow 0$ so that we get some “approximate uniformity”; in particular, all codewords in $\text{Dec}^{-1}(\mathbf{m}) \cap A_{\delta(n)}$ will have approximately the same probability mass under $\text{BSC}_p(\mathbf{y})$. By [Theorem 1.2](#) (or its proof), we can make the decay $\delta(n) \rightarrow 0$ sufficiently slow while maintaining

$$\Pr [\text{BSC}_p(\mathbf{y}) \in A_{\delta(n)}] \geq 1 - \frac{\epsilon}{2}.$$

This will ensure that we do not incur too much loss when we pass from $\text{BSC}_p(\mathbf{y})$ to the conditional distribution $\text{BSC}_p(\mathbf{y}) \mid A_{\delta(n)}$. Let us now make this formal.

Applying the Law of Total Probability, we have

$$\begin{aligned} \epsilon &\leq \Pr [\text{BSC}_p(\mathbf{y}) \in \text{Dec}^{-1}(\mathbf{m})] \\ &= \Pr [\text{BSC}_p(\mathbf{y}) \in \text{Dec}^{-1}(\mathbf{m}) \cap A_{\delta(n)}] \\ &\quad + \Pr [\text{BSC}_p(\mathbf{y}) \in \text{Dec}^{-1}(\mathbf{m}) \mid \text{BSC}_p(\mathbf{y}) \notin A_{\delta(n)}] \cdot \Pr [\text{BSC}_p(\mathbf{y}) \notin A_{\delta(n)}] \\ &\leq \Pr [\text{BSC}_p(\mathbf{y}) \in \text{Dec}^{-1}(\mathbf{m}) \mid \text{BSC}_p(\mathbf{y}) \in A] + \frac{\epsilon}{2}. \end{aligned} \quad (\text{Using } \text{Theorem 1.2})$$

Hence, we obtain on the one hand

$$\Pr [\text{BSC}_p(\mathbf{y}) \in \text{Dec}^{-1}(\mathbf{m}) \cap A_{\delta(n)}] \geq \frac{\epsilon}{2} = \Omega(1).$$

On the other hand,

$$\begin{aligned} \Pr [\text{BSC}_p(\mathbf{y}) \in \text{Dec}^{-1}(\mathbf{m}) \cap A_{\delta(n)}] &= \sum_{\mathbf{x} \in \text{Dec}^{-1}(\mathbf{m}) \cap A_{\delta(n)}} \Pr [\text{BSC}_p(\mathbf{y}) = \mathbf{x}] \\ &\leq p^{(p-\delta(n))n} (1-p)^{(1-p+\delta(n))n} \cdot |\text{Dec}^{-1}(\mathbf{m}) \cap A_{\delta(n)}| \\ &= 2^{-(H(p)-o(1)) \cdot n} \cdot |\text{Dec}^{-1}(\mathbf{m}) \cap A_{\delta(n)}|, \end{aligned}$$

where the $o(1)$ is due to the fact that $\delta(n) \rightarrow 0$ as $n \rightarrow \infty$. Rearranging then yields

$$|\text{Dec}^{-1}(\mathbf{m})| \geq |\text{Dec}^{-1}(\mathbf{m}) \cap A_{\delta(n)}| \geq \frac{\epsilon}{2} \cdot 2^{(H(p)-o(1)) \cdot n} \geq 2^{(H(p)-o(1)) \cdot n},$$

where we may absorb $\frac{\epsilon}{2}$ into $2^{o(n)}$. □

3 Proof of [Theorem 1.1](#): The Case $r < 1 - H(p)$

Our task is to design an error-correcting code. This isn’t easy to do by hand, so let’s use the probabilistic method. We will let $\text{Enc} : \{0, 1\}^k \rightarrow \{0, 1\}^n$ be a randomly chosen function, and let $\text{Dec} : \{0, 1\}^n \rightarrow \{0, 1\}^k$ be the *maximum likelihood decoder*:

$$\text{Dec}(\mathbf{x}) \stackrel{\text{def}}{=} \arg \min_{\mathbf{m} \in \{0, 1\}^k} d_H(\mathbf{x}, \text{Enc}(\mathbf{m})).$$

In the case where multiple messages \mathbf{m} attain the minimum, we select one arbitrarily. This decoder is by no means efficient, but it is a natural one to use and can be computed by brute force. Note that since $\text{Enc} : \{0, 1\}^k \rightarrow \{0, 1\}^n$ is a uniformly random function, meaning $\text{Enc}(\mathbf{m}) \sim \text{Unif}\{0, 1\}^n$ independently over all $\mathbf{m} \in \{0, 1\}^k$, there is a (small) collision probability we’ll have to address in the overall error probability of decoding.

Ultimately, our goal is to guarantee that the success probability is at least $1 - o(1)$ uniformly for all $\mathbf{m} \in \{0, 1\}^k$. As a stepping stone, we begin by establishing an average-case guarantee.

Proposition 3.1. *Suppose $r < 1 - H(p)$ and $k = rn$. Then there exists $\eta(n)$ decaying to 0 as $n \rightarrow \infty$ such that the following holds: For a uniformly random function $\text{Enc} : \{0, 1\}^k \rightarrow \{0, 1\}^n$ paired with the maximum likelihood decoder $\text{Dec} : \{0, 1\}^n \rightarrow \{0, 1\}^k$,*

$$\mathbb{E} [\Pr [\text{Dec}(\text{BSC}_p(\text{Enc}(\mathbf{m}))) \neq \mathbf{m}] \leq \eta(n), \quad \forall \mathbf{m} \in \{0, 1\}^k,$$

where the inner probability is with respect to the binary symmetric channel, and the outer expectation is with respect to the random choice of encoder.

Proof. Because we're using maximum likelihood decoding and because our encoder is chosen uniformly at random, the hope is that $\text{BSC}_p(\text{Enc}(\mathbf{m}))$ will be close to $\text{Enc}(\mathbf{m})$ and far from $\text{Enc}(\mathbf{m}')$ for $\mathbf{m}' \neq \mathbf{m}$. In particular, for a parameter $t > 0$ to be determined later, observe that $\text{Dec}(\text{BSC}_p(\text{Enc}(\mathbf{m}))) = \mathbf{m}$ holds if simultaneously $d_H(\text{BSC}_p(\text{Enc}(\mathbf{m})), \text{Enc}(\mathbf{m})) \leq t$, while $d_H(\text{BSC}_p(\text{Enc}(\mathbf{m})), \text{Enc}(\mathbf{m}')) > t$ for all $\mathbf{m}' \neq \mathbf{m}$. Taking the contrapositive and applying the Union Bound, we have that

$$\begin{aligned} & \Pr[\text{Dec}(\text{BSC}_p(\text{Enc}(\mathbf{m}))) \neq \mathbf{m}] \\ & \leq \Pr[d_H(\text{BSC}_p(\text{Enc}(\mathbf{m})), \text{Enc}(\mathbf{m})) > t] + \sum_{\mathbf{m}' \neq \mathbf{m}} \Pr[d_H(\text{BSC}_p(\text{Enc}(\mathbf{m})), \text{Enc}(\mathbf{m}')) \leq t]. \end{aligned}$$

Now choose $\delta(n), \epsilon(n) \rightarrow 0$ and set $t = (p + \delta(n)) \cdot n$ so that

$$\Pr[d_H(\text{BSC}_p(\text{Enc}(\mathbf{m})), \text{Enc}(\mathbf{m})) > t] \leq \epsilon(n),$$

which is possible by [Theorem 1.2](#) (or its proof). For the remaining terms, fix $\mathbf{m}' \neq \mathbf{m}$. By writing the probability as an expectation of an indicator random variable, and then exchanging the order of the expectations, we have

$$\mathbb{E}_{\text{Enc}} \left[\Pr_{\text{BSC}_p} [d_H(\text{BSC}_p(\text{Enc}(\mathbf{m})), \text{Enc}(\mathbf{m}')) \leq t] \right] = \mathbb{E}_{\text{BSC}_p} \left[\Pr_{\text{Enc}} [d_H(\text{BSC}_p(\text{Enc}(\mathbf{m})), \text{Enc}(\mathbf{m}')) \leq t] \right].$$

The inner probability in the right-hand side is much easier to estimate, since $\text{BSC}_p(\text{Enc}(\mathbf{m}))$ can be treated as some arbitrarily fixed codeword $\mathbf{y} \in \{0, 1\}^n$, and the randomness is taken over a uniformly random $\mathbf{x} = \text{Enc}(\mathbf{m}')$. More precisely, it is equal to the probability that a uniformly random bit-string $\text{Enc}(\mathbf{m}') \sim \text{Unif}\{0, 1\}^n$ lands in the radius- t Hamming ball around the point $\mathbf{y} = \text{BSC}_p(\text{Enc}(\mathbf{m}))$. This probability is the same regardless of what \mathbf{y} is, and is upper bounded by

$$\frac{1}{2^n} \sum_{k=0}^{(p+\delta(n)) \cdot n} \binom{n}{k} \leq 2^{(H(p)-1+\delta(n)) \cdot n}. \quad (\text{Using [Lemma 1.3](#)})$$

Combining these inequalities, we obtain

$$\mathbb{E}[\Pr[\text{Dec}(\text{BSC}_p(\text{Enc}(\mathbf{m}))) \neq \mathbf{m}]] \leq \epsilon(n) + 2^{(H(p)-1+\delta(n)) \cdot n} \stackrel{\text{def}}{=} \eta(n) \leq o(1),$$

using that $p < 1/2$ so that $H(p) < 1$ is constant independent of n . \square

To complete the proof of [Theorem 1.1](#) in the case $r < 1 - H(p)$, observe that by averaging [Proposition 3.1](#) over $\mathbf{m} \in \{0, 1\}^k$, we have

$$\mathbb{E}_{\text{Enc}} [\mathbb{E}_{\mathbf{m} \sim \text{Unif}\{0, 1\}^k} [\Pr[\text{Dec}(\text{BSC}_p(\text{Enc}(\mathbf{m}))) \neq \mathbf{m}]]] \leq \eta(n).$$

Hence, there exists an encoder $\text{Enc} : \{0, 1\}^k \rightarrow \{0, 1\}^n$ possessing the average-case guarantee

$$\mathbb{E}_{\mathbf{m} \sim \text{Unif}\{0, 1\}^k} [\Pr[\text{Dec}(\text{BSC}_p(\text{Enc}(\mathbf{m}))) \neq \mathbf{m}]] \leq \eta(n).$$

Recall that ultimately, we want this bound to hold for all $\mathbf{m} \in \{0, 1\}^k$, not in expectation.

Here's the final idea. Observe that the above bound on the expectation implies that there exists a subset of "good messages" $\mathcal{G} \subseteq \{0, 1\}^k$ with $|\mathcal{G}| = 2^{k-1}$ such that for all $\mathbf{m} \in \mathcal{G}$, we have

$$\Pr[\text{Dec}(\text{BSC}_p(\text{Enc}(\mathbf{m}))) \neq \mathbf{m}] \leq 2 \cdot \eta(n).$$

Indeed, if more than half of the elements of $\{0, 1\}^k$ have failure probability exceeding $2 \cdot \eta(n)$, then we would obtain a contradiction with the aforementioned bound on the expectation. Hence, we have found an encoder with the desired properties for a set of messages $\mathcal{G} \subseteq \{0, 1\}^k$ of size 2^{k-1} . By mapping \mathcal{G} bijectively to $\{0, 1\}^{k-1}$, it follows we have constructed a $[n, k-1]$ -code with worst-case failure probability $2 \cdot \eta(n) \leq o(1)$. Replacing k with $k-1$, which incurs $o(1)$ in the rate r , then completes the proof.