

Lecture 3: Shannon's Theorem

October 9, 2006

Lecturer: Venkatesan Guruswami

Scribe: Widad Machmouchi

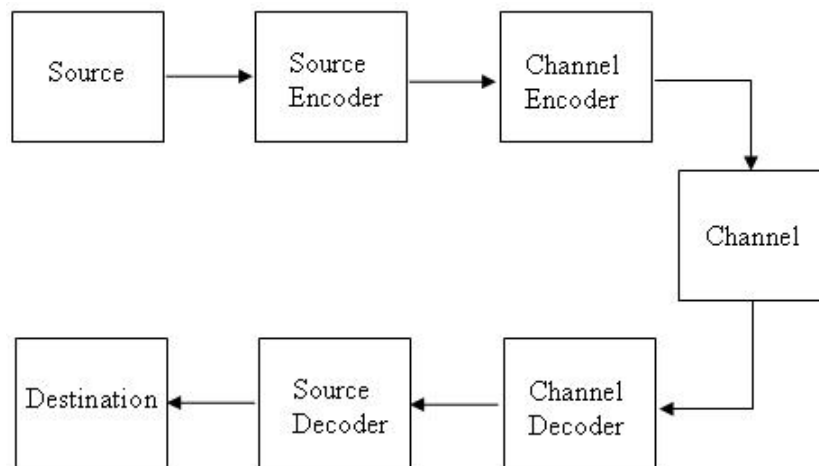
1 Communication Model

The communication model we are using consists of a source that generates digital information. This information is sent to a destination through a channel. The communication can happen in the time domain or space domain.

The channel can be associated with noise. So we have to cases :

- **Noiseless case** The channel in this case transmits symbols without causing any errors. One would need to exploit the redundancy in the source to economize the length of the transmission. This is done through data compression, also called source coding. The information is decompressed at destination.
- **Noisy case** The channel in this case introduces noise that causes errors in the received symbols at the destination. To reduce the errors incurred due to noise, one should add systematic redundancy to the information to be sent. This is done through channel coding.

The following diagram shows the modules of the communication model:



We Will focus on the channel coding problem, assuming message to be communicated is just a collection of symbols over a finite alphabet. We can communicate using the above two-stage

scheme and treat source and channel coding in isolation due to the Source-Channel Coding theorem.

Source-Channel Coding Theorem

For a source with entropy no greater than the capacity of the channel, dividing the transmission process to source coding and then channel coding is optimal and will achieve a probability of error tending to zero for a large block length.

Then the part of the previous scheme we'll be considering is:



Shannon put forth a stochastic model of the channel. For us, it suffices to talk about discrete memoryless channels. Such channels have an input alphabet \mathcal{X} , an output alphabet \mathcal{Y} and a probability transition matrix describing the output distribution for every input. Each symbol is sent over the channel independently of the previous symbols sent. Thus the channel is prescribed by a $|\mathcal{X}| \times |\mathcal{Y}|$ stochastic matrix where each row sums to 1.

Probability transition matrix:

$$|\mathcal{X}| \left\{ \overbrace{\left(\begin{matrix} p(y/x) \end{matrix} \right)}^{|\mathcal{Y}|} \right.$$

2 Examples of Channels

Channels are often described by input-output diagrams.

2.1 Binary Symmetric Channel (BSC)

the BSC takes as input one bit (0/1) and flips it with probability p , $0 \leq p \leq \frac{1}{2}$. p is called the crossover probability; we write BSC_p .

2.2 Binary Erasure Channel (BEC)

The BEC takes as input one bit (0/1) and flips it to ? with probability ε , $0 \leq \varepsilon \leq \frac{1}{2}$. ε is called the erasure probability; we write BEC_ε .

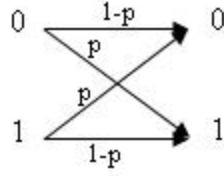


Figure 1: Diagram for BSC_p

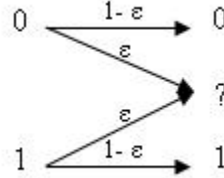


Figure 2: Diagram for BEC_ϵ

2.3 Continuous Output Channel

The continuous output channel takes as input a symbol from a finite alphabet and maps it, according to a specific noise distribution, to a real number. One example is the Binary Input Additive White Gaussian Noise (BIAWGN) channel. The noise has a normal distribution.

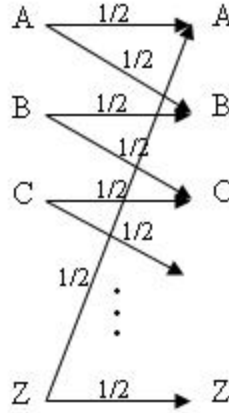
$$\begin{array}{ccccc} \Sigma = \{-1, 1\} & \rightarrow & \boxed{\text{Channel}} & \rightarrow & \mathbb{R} \\ x & \longrightarrow & & & y = x + z; z \in N(0, \sigma^2) \end{array}$$

Hence the probability density function of the noise is given by :

$$\Pr(z) = \Pr(y|x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y-x)^2}{2\sigma^2}\right)$$

2.4 Noisy Typewriter Channel

The noisy typewriter channel is given by the following diagram:



Zero-error communication is possible if we just send A, C, E, ..., Y or B, D, F, ..., Z. So we will end up with only 13 possibilities for the sent symbols. Thus the capacity of the channel is at least $\log_2 13$ bits. In fact, one can prove that this rate is the maximum possible and the capacity of the channel is exactly $\log_2 13$ bits.

Zero-error communication at a positive rate is not possible with BSC_p , since for every pair of strings x and y , there is a positive probability that x gets distorted into y . This probability of miscommunication can be reduced arbitrarily by high enough order repetition code.

Consider mapping 0 to m zeroes and 1 to m ones. At the destination, decoding is done by majority, i.e. if the number of received 0's is greater than $\frac{m}{2}$, we decide on a 0, else we decide on a 1. Hence, the probability of error is given by: $\Pr(\text{error}) = \sum_{i=\frac{m}{2}}^m \binom{m}{i} p^i (1-p)^{m-i}$. This probability tends to 0 as m tends to ∞ but this causes the rate to tend to zero!!

Can we achieve any desired probability of error while maintaining a positive rate? The answer is yes. In fact, the largest possible rate was precisely described in Shannon's work.

3 Shannon Capacity Theorem

We will start by defining the binary entropy function.

Definition 3.1. For $0 \leq x \leq 1$, the entropy binary function, denoted $H(x)$ is given by

$$H(x) = x \log_2 \frac{1}{x} + (1-x) \log_2 \frac{1}{1-x}.$$

Theorem 3.2 (Shannon Capacity Theorem). For every $p, 0 \leq p \leq \frac{1}{2}, \epsilon > 0$, there exists $\delta > 0$ such that for all large n , there exist an encoding function $E : \{0, 1\}^k \rightarrow \{0, 1\}^n$ and a decoding function $D : \{0, 1\}^n \rightarrow \{0, 1\}^k$ for $k = (1 - H(p + \epsilon))n$ such that $\forall m \in \{0, 1\}^k$,

$$\Pr_{\text{noise of the } BSC_p} (D(E(m) + \text{noise}) \neq m) \leq 2^{-\delta n}.$$

We will first prove the converse of Shannon theorem to give an intuition why $1 - H(p)$ is the best one should hope for. For this purpose we will need the following lemma.

Lemma 3.3. For $0 \leq p \leq \frac{1}{2}$, $\sum_{i=0}^{pn} \binom{n}{i} \leq 2^{H(p)n}$.

Proof

$$\begin{aligned}
1 &= (p + (1 - p))^n \\
&\geq \sum_{i=0}^{pn} \binom{n}{i} p^i (1 - p)^{n-i} \\
&= \sum_{i=0}^{pn} \binom{n}{i} p^i (1 - p)^n \left(\frac{p}{1 - p} \right)^i \\
&\geq \sum_{i=0}^{pn} \binom{n}{i} p^i (1 - p)^n \left(\frac{p}{1 - p} \right)^{pn} \\
&= \sum_{i=0}^{pn} \binom{n}{i} 2^{-H(p)n}
\end{aligned}$$

Note $2^{H(p)n} = 2^{(-p \log p - (1-p) \log (1-p))n} = \frac{1}{p^{pn} (1-p)^{(1-p)n}}$

We will also need the Chernoff bound.

Chernoff bound If X_1, X_2, \dots, X_n are i.i.d 0/1 random variables with $\Pr[X_i = 1] = p$, then $\forall \varepsilon > 0$

$$\begin{aligned}
\Pr\left[\sum_{i=1}^n X_i \geq (p + \varepsilon)n\right] &\leq 2^{-\frac{\varepsilon^2 n}{4p}} \\
\Pr\left[\sum_{i=1}^n X_i \leq (p - \varepsilon)n\right] &\leq 2^{-\frac{\varepsilon^2 n}{2p}}
\end{aligned}$$

Proof of the converse of Shannon theorem

By the Chernoff bound, with high probability, $y = E(m) + \text{noise}$ will lie in a shell of radii $((p - \varepsilon)n, (p + \varepsilon)n)$ around $E(m)$. The number of n -bit strings lying within a distance pn of $E(m)$ is the volume of the hamming ball $B(0, pn)$.

$\text{Vol}_2(B(0, pn)) = \sum_{i=0}^{pn} \binom{n}{i}$ which is at most $2^{H(p)n}$, by the previous lemma. Then at most $2^{H(p)n}$ n -bit strings can be decoded to m with an exponentially decreasing probability of error. Hence the maximum number of codewords, 2^k , is $\frac{2^n}{2^{H(p)n}}$, implying $k \leq [(1 - H(p))n]$.

To prove Shannon theorem, we need to establish the following:

$$\text{Vol}_2(B(0, pn)) \approx 2^{H(p)n}$$

We know from a previous lemma that $\text{Vol}_2(B(0, pn)) \leq 2^{H(p)n}$. Add to this the fact $\binom{n}{pn} \geq 2^{H(p)n - o(n)}$. This fact follows from applying Stirling approximation of $n!$ given by:

$$n! = \sqrt{2\pi n} \left(\frac{n}{e}\right)^n \left(1 + \Theta\left(\frac{1}{n}\right)\right)$$

We will continue the proof of Shannon theorem next lecture.