

x Lecture 6, Oct. 7, 2003

Theorem: The Channel Coding Theorem:

All rates below capacity C are achievable, i.e., for every $R < C$ and any $\epsilon > 0$, there exists a sequence of $(2^{nR}, n)$ codes with maximum probability of error $\lambda^{(n)} \rightarrow 0$.

Theorem: Converse to the channel coding theorem
~~Conversely~~, Any sequence of $(2^{nR}, n)$ codes with $\lambda^{(n)} \rightarrow 0$ must have ~~$R < C$~~ $R \leq C$.

Proof of Coding ~~theorem~~ the achievability

Fix $p(x)$ and generate a $(2^{nR}, n)$ code at random according to the distribution, \mathcal{C}

$$p(x^n) = \prod_{i=1}^n p(x_i)$$

$$\mathcal{C} = \begin{bmatrix} x_1(1) & x_2(1) & \dots & x_n(1) \\ \vdots & \vdots & & \vdots \\ x_1(2^{nR}) & x_2(2^{nR}) & & x_n(2^{nR}) \end{bmatrix}$$

Each entry in the matrix is generated i.i.d.

according to $p(x)$. nR

$$P_r(\mathcal{C}) = \prod_{w=1}^{2^{nR}} \prod_{i=1}^n p(x_i(w))$$

This would be the sequence of events in the channel coding process:

- 1) A random code \mathcal{C} is generated at random according to $p(x)$
- 2) The code \mathcal{C} is then revealed to both the sender and receiver. Both sender and receiver know the transition matrix $P(y|x)$.
- 3) A message w is chosen according to a uniform distribution:

$$P_r(W=w) = 2^{-nR} \quad w=1, 2, \dots, 2^{nR}$$

- 4) The w th codeword $X^n(w)$, corresponding to the w th row of \mathcal{C} , is sent over the channel.
- 5) The receiver receives a sequence Y^n according

to

$$P(Y^n | X^n(w)) = \prod_{i=1}^n P(y_i | x_i(w))$$

- 6) The receiver makes a decision about the message sent. In practice, the receiver uses maximum likelihood ~~decoding~~ ^{which is optimal} decoding. However, to make the analysis simple, we assume that the receiver uses typical set decoding. We describe this as:

The receiver declares that the index \hat{W} was sent if the following conditions are satisfied:

- $(X^n(\hat{W}), Y^n)$ is jointly typical.
- There is no other index k such that $(X^n(k), Y^n) \in A_{\epsilon}^{(n)}$

If no such \hat{W} exist or if there are more than one such \hat{W} , then an error is declared.

1) There is a decoding error if $\hat{W} \neq W$. Denote by \mathcal{E} the event $\hat{W} \neq W$.

Analysis of the probability of error:

$$\begin{aligned}
 P_r(\mathcal{E}) &= \sum_{\mathcal{C}} P(\mathcal{C}) P_e^{(n)}(\mathcal{C}) \\
 &= \sum_{\mathcal{C}} P(\mathcal{C}) \frac{1}{2^{nR}} \sum_{w=1}^{2^{nR}} \lambda_w(\mathcal{C}) \\
 &= \frac{1}{2^{nR}} \sum_{w=1}^{2^{nR}} \sum_{\mathcal{C}} P(\mathcal{C}) \lambda_w(\mathcal{C})
 \end{aligned}$$

where $P_e^{(n)}(\mathcal{C})$ is defined for typical set decoding.

The way we have constructed the code results in symmetry w.r.t. Codewords, i.e., the average probability of error (averaged over all codes) does not depend on the particular index that was sent, i.e., $\sum_C P(C) \lambda_w(C)$ does not depend on w . ~~So~~, Thus, we can assume that $W=1$ was sent. Then,

$$P_r(E) = \frac{1}{2^{nR}} \sum_{w=1}^{2^{nR}} \sum_C P(C) \lambda_w(C)$$

$$= \sum_C P(C) \lambda_1(C)$$

$$= P_r(E|W=1)$$

Define the event E_i as:

$$E_i = \{(x^n(i), y^n) \in A_E^{(n)}\} \quad i \in \{1, 2, \dots, 2^{nR}\}$$

E_i is the event that the i th codeword and y^n are jointly typical.

Then ~~an error~~ a decoding error occurs either when E_1^c occurs (when the transmitted and received sequences are not jointly typical) or $E_2 \cup E_3 \dots \cup E_{2^{nR}}$ occurs, i.e., when another

Codeword is jointly typical with y^n .

Hence, letting $P(\mathcal{E}) = P_r(\mathcal{E}|W=1)$, we have

$$\begin{aligned} P_r(\mathcal{E}|W=1) &= P(E_1^c \cup E_2 \cup E_3 \dots \cup E_{2^{nR}}) \\ &\leq P(E_1^c) + \sum_{i=2}^{2^{nR}} P(E_i) \end{aligned}$$

From the joint AEP, $P(\bar{E}_1^c) \rightarrow 0$, hence

$$P(E_1^c) \leq \epsilon \quad \text{for } n \text{ sufficiently large.}$$

Since by code construction $X^n(1)$ and $X^n(i)$ are independent so are y^n and $X^n(i)$, $i \neq 1$.

Hence, the probability that $X^n(i)$ and y^n are jointly typical is $\leq 2^{-n(I(X;Y) - 3\epsilon)}$

$$\begin{aligned} P(\mathcal{E}) &= P(\mathcal{E}|W=1) \leq P(E_1^c) + \sum_{i=2}^{2^{nR}} P(E_i) \\ &\leq \epsilon + \sum_{i=2}^{2^{nR}} 2^{-n(I(X;Y) - 3\epsilon)} \\ &= \epsilon + (2^{nR} - 1) 2^{-n(I(X;Y) - 3\epsilon)} \\ &\leq \epsilon + \frac{2^{nR}}{2} \cdot 2^{-n(I(X;Y) - R)} \end{aligned}$$

if n is sufficiently large and $R < I(X;Y) - 3\epsilon$ then $P(\mathcal{E}) \leq 2\epsilon$.

Hence if $R < I(X;Y)$, we can choose ϵ and n such that the average probability of error is

less than 2ϵ .

Now we finish the proof by:

1) Choosing $p(x)$ to be $p^*(x)$, i.e., the input probability distribution that maximizes $I(X;Y)$ (to get $C = I(X;Y) |_{p^*(x)}$). Then, we can replace $R < I(X;Y)$ by $R < C$.

2) Since probability of error averaged over all codebooks is small, there should be one that results in small probability of error.

"We can find" such a good codebook^{C*} by exhaustive search^{over} of all $(2^{nR}, n)$ codes.

This way, we get rid of average over codebooks.

3) ~~Throw~~ Take the best codebook^(C*). For this codebook^V, we have

$$\frac{1}{2^{nR}} \sum_i \lambda_i(C^*) \leq 2\epsilon$$

through throw away ~~the~~ half of the ~~codebook~~ codewords of C^* (the worst ones). The

remaining half have per conditional probability of error λ_i (associated to $X^n(i)$) less than 4ϵ .

If this was not true, these codewords themselves could contribute 2ϵ or more to the sum.

If we re-index the remaining codewords, we get 2^{nR-1} codewords. So, we get a codebook of rate $R - \frac{1}{n} \xrightarrow{n \rightarrow \infty} R$ with maximal error probability $\lambda^{(n)} \leq 4\epsilon$.

~~~~~  
We have proved the achievability, i.e., the possibility of error free transmission for  $R < C$ . Now, we prove that error-free transmission is not possible if  $R > C$ . This is the converse to coding theorem.  
Theorem: Any sequence of  $(2^{nR}, n)$  codes with  $\lambda^{(n)} \rightarrow 0$  must have  $R \leq C$ .

~~~~~  
In order to prove the converse channel coding theorem, we use the following two lemmas.

Lemma: (Fano's inequality): For a discrete memory less channel with a codebook C and the input message uniformly distributed, let $P_e^{(n)} = P_r(W \neq g(Y^n))$. Then:

$$H(X^n | Y^n) \leq 1 + P_e^{(n)} nR$$

Define the random variable E as

$$E = \begin{cases} 1 & \text{if } \hat{w} = g(y^n) \neq w \\ 0 & \text{if } \hat{w} = w \end{cases}$$

Then,

$$\begin{aligned} H(E, W | Y^n) &= H(W | Y^n) + H(E | W, Y^n) \\ &= H(E | Y^n) + H(W | E, Y^n) \end{aligned}$$

$H(E | W, Y^n) = 0$ since E is a function of W and $g(y^n)$.

So,

$$H(W | Y^n) = H(E | Y^n) + H(W | E, Y^n)$$

$H(E | Y^n) \leq 1$ since E is a binary r.v.

$$\begin{aligned} H(W | E, Y^n) &= P(E=0) H(W | E=0, Y^n) \\ &\quad + P(E=1) H(W | E=1, Y^n) \\ &= (1 - P_e^{(n)}) \times 0 + P_e^{(n)} H(W | E=1, Y^n) \\ &\leq P_e^{(n)} \frac{1}{2} \log(|W| - 1) \\ &\leq P_e^{(n)} n R \end{aligned}$$

So, $H(W | Y^n) \leq 1 + P_e^{(n)} n R$

Since for a given fixed code $X^n(w)$ is a

function of W ,

$$H(X^n(W) | Y^n) \leq H(W | Y^n)$$

So,

$$H(X^n | Y^n) \leq 1 + P_e^{(n)} n R$$

Lemma: Let Y^n be the result of passing X^n through a DMC. Then

$$I(X^n; Y^n) \leq nC \quad \text{for all } P(X^n).$$

Interpretation: This means that by using a DMC several times, we cannot increase the capacity per use.

Proof:

$$\begin{aligned} I(X^n; Y^n) &= H(Y^n) - H(Y^n | X^n) \\ &= H(Y^n) - \sum_{i=1}^n H(Y_i | Y_1, \dots, Y_{i-1}, X^n) \\ &= H(Y^n) - \sum_{i=1}^n H(Y_i | X_i) \end{aligned}$$

{ since the entropy of a collect of r.v.'s is less than the sum of their entropies } \rightarrow

$$\begin{aligned} &\leq \sum_{i=1}^n H(Y_i) - \sum_{i=1}^n H(Y_i | X_i) \\ &= \sum_{i=1}^n I(X_i; Y_i) \leq nC. \end{aligned}$$

Now, we proceed with the proof of the converse to Channel Coding theorem, i.e., we show that any sequence $(2^{nR}, n)$ of codes with $\lambda^{(n)} \rightarrow 0$ must satisfy $R \leq C$.

Proof: if $\lambda^{(n)} \rightarrow 0$ then so does the average probability of error $P_e^{(n)} = \frac{1}{M} \sum_{i=1}^M \lambda_i$.

For each n , let W be drawn uniformly from $\{1, 2, \dots, 2^{nR}\}$. Then $P_e^{(n)} = P_r(\hat{W} \neq W)$.

Hence,

$$nR = H(W) = H(W|Y^n) + I(W; Y^n)$$

$$\leq 1 + P_e^{(n)} nR + I(X^n(W); Y^n)$$

$$\leq 1 + P_e^{(n)} nR + nC$$

Dividing both sides by n ,

$$R \leq P_e^{(n)} R + \frac{1}{n} + C \quad (A)$$

letting $n \rightarrow \infty$ the first two terms $\rightarrow 0$ and we get $R \leq C$.

We can write (A) as

$$P_e^{(n)} \geq 1 - \frac{C}{R} - \frac{1}{nR}$$

So, if $R > C$ then $P_e^{(n)}$ is bounded away from zero for sufficiently large n , and, hence for all n . The latter comes from the fact that if $P_e^{(n)} \rightarrow 0$ for small n , we could construct long codes by concatenating short codes.

X Lecture 7, Oct. 14, 2003

Continuous Sources

A continuous random variable X is characterised by a Cumulative Distribution Function $F_X(x)$,

$$F_X(x) = P_r(X \leq x)$$

the derivative of $F_X(x)$ w.r.t. x is called the probability density function (pdf),

$$f_X(x) = \frac{\partial F_X(x)}{\partial x}$$

when there is no ambiguity, we drop the subscript and denote the pdf by $f(x)$.

It is clear that,

$$F(x) = \int_{-\infty}^x f(x) dx$$

therefore

$$\int_S f(x) dx = 1$$