

# Information Theory and Networks

## Lecture 24: Channel Capacity

Matthew Roughan

[<matthew.roughan@adelaide.edu.au>](mailto:matthew.roughan@adelaide.edu.au)

[http://www.maths.adelaide.edu.au/matthew.roughan/  
Lecture\\_notes/InformationTheory/](http://www.maths.adelaide.edu.au/matthew.roughan/Lecture_notes/InformationTheory/)

School of Mathematical Sciences,  
University of Adelaide

October 9, 2013

# Part I

## Channel Capacity

To make no mistakes is not in the power of man; but from their errors and mistakes the wise and good learn wisdom for the future.

*Plutarch*

# Section 1

## Channel Coding Theorem

# Capacity

## Definition (Operational Channel Capacity)

The highest rate of bits we can send per input symbol, with an arbitrarily low probability of error is called the **operational channel capacity**.

## Definition (Information Capacity)

The **information capacity** of a discrete memoryless channel with inputs  $X \in \mathcal{X}$  and outputs  $Y \in \mathcal{Y}$ , and channel transition matrix  $p(Y|X)$  is

$$C = \max_{p_X(x)} I(X; Y)$$

where  $I(X; Y)$  is the mutual information of  $X$  and  $Y$ .

Capacity

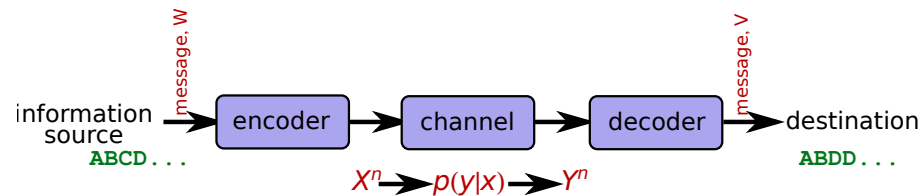
Definition (Operational Channel Capacity)  
The highest rate of bits we can send per input symbol, with an arbitrarily low probability of error is called the operational channel capacity.

Definition (Information Capacity)  
The information capacity of a discrete memoryless channel with inputs  $X \in \mathcal{X}$  and outputs  $Y \in \mathcal{Y}$ , and channel transition matrix  $p(Y|X)$  is

$$C = \max_{p_X(x)} I(X; Y)$$

where  $I(X; Y)$  is the mutual information of  $X$  and  $Y$ .

# Digital Communications Channels



## Definition (Discrete Channel)

A **discrete channel** is a system with an input alphabet  $\mathcal{X}$ , and output alphabet  $\mathcal{Y}$ , and a probability transition matrix  $p(y|x)$  that describes the probability of observing the output symbol  $y \in \mathcal{Y}$  given input  $x \in \mathcal{X}$ .

We denote a Discrete Memoryless Channel (DMC) by the triple  $(\mathcal{X}, p(y|x), \mathcal{Y})$ .

Digital Communications Channels

Definition (Discrete Channel)  
A discrete channel is a system with an input alphabet  $\mathcal{X}$ , and output alphabet  $\mathcal{Y}$ , and a probability transition matrix  $p(y|x)$  that describes the probability of observing the output symbol  $y \in \mathcal{Y}$  given input  $x \in \mathcal{X}$ .

We denote a Discrete Memoryless Channel (DMC) by the triple  $(\mathcal{X}, p(y|x), \mathcal{Y})$ .

# Digital Communications Channels

We will work with DMC (Discrete Memoryless Channels) with no feedback  $(\mathcal{X}, p(y|x), \mathcal{Y})$ . Then

## Definition

The  $n$ th extension of a DMC is the channel  $(\mathcal{X}^n, p(y^{(n)}|x^{(n)}), \mathcal{Y}^n)$  where

$$p(y_k|x^{(k)}, y^{(k-1)}) = p(y_k|x_k), \text{ for } k = 1, 2, \dots, n$$

and/or

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

Digital Communications Channels

We will work with DMC (Discrete Memoryless Channels) with no feedback  $(\mathcal{X}, p(y|x), \mathcal{Y})$ . Then

**Definition**

The  $n$ th extension of a DMC is the channel  $(\mathcal{X}^n, p(y^{(n)}|x^{(n)}), \mathcal{Y}^n)$  where

$$p(y_k|x^{(k)}, y^{(k-1)}) = p(y_k|x_k), \text{ for } k = 1, 2, \dots, n$$

and/or

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

There are extensions of the theory to channels with memory, but these get scary as we then have to define a model for memory, and avoid any pathological cases (such as a panic button, where one input can kill the channel – I've heard of at least one case where a router's line card would die given a certain packet as input).

The extension is basically just the channel, where we use block codewords of length  $n$ .

# Channel Codes

## Definition (Channel Code)

A  $(M, n)$  code for channel  $(\mathcal{X}, p(y|x), \mathcal{Y})$  consists of

- 1 An index set  $\{1, 2, \dots, M\}$
- 2 An encoding function with block size  $n$

$$X^n : \{1, 2, \dots, M\} \rightarrow \mathcal{X}^n$$

yielding codewords  $\{X^n(1), X^n(2), \dots, X^n(M)\}$ , called the **codebook**.

- 3 A decoding function

$$g : \mathcal{Y}^n \rightarrow \{1, 2, \dots, M\}$$

which is a deterministic rule which assigns a guess to each possible received vector.

Channel Codes

**Definition (Channel Code)**

A  $(M, n)$  code for channel  $(\mathcal{X}, p(y|x), \mathcal{Y})$  consists of

- 1 An index set  $\{1, 2, \dots, M\}$
- 2 An encoding function with block size  $n$ 

$$X^n : \{1, 2, \dots, M\} \rightarrow \mathcal{X}^n$$
- 3 yielding codewords  $\{X^n(1), X^n(2), \dots, X^n(M)\}$ , called the **codebook**.
- 4 A decoding function
 
$$g : \mathcal{Y}^n \rightarrow \{1, 2, \dots, M\}$$

which is a deterministic rule which assigns a guess to each possible received vector.

1. The index set is just an abstraction of the input symbols – we are assuming that there are  $M$  such.
2. Each input symbol will be translated into a block of  $n$  symbols from  $\mathcal{X}$  to be transmitted over the channel.
3. The decoding function estimates which input codeword was most likely from the block  $y^{(n)}$  of output symbols.

For example, the input symbols might be ASCII, in which case  $M = 256$ , and then these are translated into the standard 8-bit binary codewords, i.e.,  $\mathcal{X}^n = \{0, 1\}^8$ .

It is important to note that the set of codewords that we use might not be the full set of possible codewords from  $\mathcal{X}^n$ , i.e., we could choose  $n > \log_2(M)$ . But  $g(\cdot)$  must contain mappings from all of  $\mathcal{Y}^n$ .

# Errors

## Definition

The conditional probability of error given that index  $i$  is sent is

$$\lambda_i = P(g(Y^n) \neq i \mid X^n = X^n(i)) = \sum_{y^n} p(y^n | x^n(i)) I(g(y^n) \neq i)$$

where  $I(\cdot)$  is an indicator function.

An indicator function for even  $A$  is

$$I(A) = \begin{cases} 1, & \text{if } A, \\ 0, & \text{otherwise} \end{cases}$$

# Errors

## Definition

The maximal probability of error  $\lambda^{(n)}$  for an  $(M, n)$  code is defined as

$$\lambda^{(n)} = \max_{i \in \{1, 2, \dots, M\}} \lambda_i$$

and the average probability of error  $P_e^{(n)}$  is

$$P_e^{(n)} = \frac{1}{M} \sum_{i=1}^M \lambda_i = P(I \neq g(Y^n))$$

where  $I$  is a random index uniformly chosen from  $\{1, 2, \dots, M\}$ .

Note that  $P_e^{(n)} \leq \lambda^{(n)}$ . We might suspect the two types of error will behave quite differently, but it turns out a small average implies a small maximum (see [CT91] for proof).

# Rate and Capacity

## Definition (Rate)

The **rate**  $R$  of an  $(M, n)$  code is

$$R = \frac{\log M}{n} \text{ bits per transmission}$$

A rate is said to be **achievable** if there exists a sequence of  $(\lceil 2^{nR} \rceil, n)$  codes such that the maximal probability of error  $\lambda^{(n)} \rightarrow 0$  as  $n \rightarrow \infty$ .

## Definition (Operational Channel Capacity)

The **capacity** of a DMC is the supremum of all the achievable rates.

Rate and Capacity

Definition (Rate)  
The rate  $R$  of an  $(M, n)$  code is  
$$R = \frac{\log M}{n} \text{ bits per transmission}$$
  
A rate is said to be **achievable** if there exists a sequence of  $(\lceil 2^{nR} \rceil, n)$  codes such that the maximal probability of error  $\lambda^{(n)} \rightarrow 0$  as  $n \rightarrow \infty$ .

Definition (Operational Channel Capacity)  
The capacity of a DMC is the supremum of all the achievable rates.

For brevity we shall write  $(2^{nR}, n)$  codes to mean  $(\lceil 2^{nR} \rceil, n)$ .

The idea of **achievable** rate is that we should be able to achieve this rate with arbitrarily low errors for large enough blocks. So this is an asymptotic definition, and in fact the blocks might have to be very large to achieve anything close to it.

The channel capacity is then the largest rate we can achieve (with arbitrarily low loss).

# Shannon's Second Theorem

## Theorem (Shannon's Channel Coding Theorem)

All rates below capacity  $C = \max_{p_X(x)} I(X; Y)$  are achievable. Specifically, for every rate  $R < C$ , there exists a sequence of  $(2^{nR}, n)$  codes with maximum probability of error  $\lambda^{(n)} \rightarrow 0$ .

Conversely, any sequence of  $(2^{nR}, n)$  codes with maximum probability of error  $\lambda^{(n)} \rightarrow 0$  must have  $R \leq C$ .

Shannon's Second Theorem

Theorem (Shannon's Channel Coding Theorem)  
All rates below capacity  $C = \max_{p_X(x)} I(X; Y)$  are achievable. Specifically for every rate  $R < C$ , there exists a sequence of  $(2^{nR}, n)$  codes with maximum probability of error  $\lambda^{(n)} \rightarrow 0$ .  
Conversely, any sequence of  $(2^{nR}, n)$  codes with maximum probability of error  $\lambda^{(n)} \rightarrow 0$  must have  $R \leq C$ .

This can also be extended so that if we can tolerate a given probability  $P_e$  of error, then we can transmit at rate:

$$R(P_e) = \frac{C}{1 - H(P_e)}$$

[Mac11, ]

# Shannon's Second Theorem

## Shannon's Channel Coding Theorem.

Full proof [CT91, pp.198-209], but some intuition follows:

- 1 We want to exploit the law of large numbers for larger blocks to obtain something like convergence to accurate estimates.
- 2 We can't increase capacity of a memoryless channel by using it multiple times, independently.
- 3 So there need to be some structure in what we send, and we are achieving this through our set of codewords.
- 4 By choosing a set of codewords that are reasonable distances apart, we hope that the errors result in sequences that are closer to the real codeword than any other.
- 5 It turns out random codewords are good enough.

Shannon's Channel Coding Theorem

Full proof [CT91, pp.198-209], but some intuition follows:

- 1 We want to exploit the law of large numbers for larger blocks to obtain something like convergence to accurate estimates.
- 2 We can't increase capacity of a memoryless channel by using it multiple times, independently.
- 3 So there need to be some structure in what we send, and we are achieving this through our set of codewords.
- 4 By choosing a set of codewords that are reasonable distances apart, we hope that the errors result in sequences that are closer to the real codeword than any other.
- 5 It turns out random codewords are good enough.

# Random Codes

- 1 Fix  $p(x)$ , and generate a random  $(2^{nR}, n)$  code by taking

$$P(X^n(i) = x_1 x_2 \dots x_n) = \prod_{k=1}^n p(x_k) \text{ for each } i \in \{1, 2, \dots, M = 2^{nR}\}$$

- 2 Write codewords as a  $2^{nR} \times n$  matrix, with IID rows

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

- 3 The probability of a particular code is

$$P(C) = \prod_{w=1}^{2^{nR}} \prod_{k=1}^n p(x_k(w))$$

Random Codes

- 1 Fix  $p(x)$ , and generate a random  $(2^{nR}, n)$  code by taking  $P(X^n(i) = x_1 x_2 \dots x_n) = \prod_{k=1}^n p(x_k)$  for each  $i \in \{1, 2, \dots, M = 2^{nR}\}$
- 2 Write codewords as a  $2^{nR} \times n$  matrix, with IID rows  $C = \begin{bmatrix} x_1(1) & x_2(1) & \dots & x_n(1) \\ x_1(2) & x_2(2) & \dots & x_n(2) \\ \dots & \dots & \dots & \vdots \\ x_1(2^{nR}) & x_2(2^{nR}) & \dots & x_n(2^{nR}) \end{bmatrix}$
- 3 The probability of a particular code is  $P(C) = \prod_{w=1}^{2^{nR}} \prod_{k=1}^n p(x_k(w))$

# Using Random Codes

To use code  $\mathcal{C}$

- 1 Assume receiver and sender both know the code, and also the transition probabilities  $p(y|x)$ .
- 2 Assume message chosen according to uniform distribution

$$P(W = w) = 2^{-nR}, \text{ for } w = 1, 2, \dots, 2^{nR}$$

and the  $w$ th codeword  $x^n(w)$  is sent.

- 3 Receiver receives  $Y^n$  according to the distribution

$$P(y^{(n)}|x^{(n)}(w)) = \prod_{i=1}^n p(y_i|x_i(w))$$

- 4 Receiver decodes by guessing that  $w$  is the input that generates a **jointly typical** sequence  $(x^{(n)}(w), y^{(n)})$ .

Using Random Codes

To use code  $\mathcal{C}$

- 1 Assume receiver and sender both know the code, and also the transition probabilities  $p(y|x)$ .
- 2 Assume message chosen according to uniform distribution  $P(W = w) = 2^{-nR}$ , for  $w = 1, 2, \dots, 2^{nR}$  and the  $w$ th codeword  $x^n(w)$  is sent.
- 3 Receiver receives  $Y^n$  according to the distribution  $P(y^{(n)}|x^{(n)}(w)) = \prod_{i=1}^n p(y_i|x_i(w))$
- 4 Receiver decodes by guessing that  $w$  is the input that generates a **jointly typical** sequence  $(x^{(n)}(w), y^{(n)})$ .

The decoding procedure isn't optimal (we should optimise to minimise errors), but is enough to demonstrate the result.

# Joint AEP (see [CT91, Theorem 8.6.1, pp.195-196])

## Definition (Jointly Typical)

The set  $A_\epsilon^{(n)}$  of **jointly typical** sequences WRT to  $p(x, y)$  is the set of sequences of  $n$  pairs  $(x_i, y_i)$  with entropies  $\epsilon$ -close to the true entropy, i.e.,

$$A_\epsilon^{(n)} = \left\{ (x^{(n)}, y^{(n)}) \mid d_X < \epsilon, d_Y < \epsilon, d_{X,Y} < \epsilon, \right\}$$

where

$$d_X = \left| -\frac{1}{n} \log p(x^{(n)}) - H(X) \right|$$

$$d_Y = \left| -\frac{1}{n} \log p(y^{(n)}) - H(Y) \right|$$

$$d_{X,Y} = \left| -\frac{1}{n} \log p(x^{(n)}, y^{(n)}) - H(X, Y) \right|$$

Joint AEP (see [CT91, Theorem 8.6.1, pp.195-196])

Definition (Jointly Typical)

The set  $A_\epsilon^{(n)}$  of jointly typical sequences WRT to  $p(x, y)$  is the set of sequences of  $n$  pairs  $(x_i, y_i)$  with entropies  $\epsilon$ -close to the true entropy, i.e.,

$$A_\epsilon^{(n)} = \left\{ (x^{(n)}, y^{(n)}) \mid d_X < \epsilon, d_Y < \epsilon, d_{X,Y} < \epsilon, \right\}$$

where

$$d_X = \left| -\frac{1}{n} \log p(x^{(n)}) - H(X) \right|$$

$$d_Y = \left| -\frac{1}{n} \log p(y^{(n)}) - H(Y) \right|$$

$$d_{X,Y} = \left| -\frac{1}{n} \log p(x^{(n)}, y^{(n)}) - H(X, Y) \right|$$

This mirrors our earlier definition of typical sequences, but now we are looking for sequences, where each component is typical (given its marginal distribution) and the combined sequence is typical given one is generated by transmission of the other over the channel.

# Joint AEP (see [CT91, Theorem 8.6.1, pp.195-196])

## Theorem (Joint AEP)

Let  $(X^{(n)}, Y^{(n)})$  be sequences of length  $n$  drawn IID according to  $p(x^{(n)}, y^{(n)}) = \prod_i p(x_i, y_i)$ , and choose  $A_\epsilon^{(n)}$  to be the set of jointly typical sequences WRT to  $p(x, y)$  then

- 1  $P((X^{(n)}, Y^{(n)}) \in A_\epsilon^{(n)}) \rightarrow 1$  as  $n \rightarrow \infty$
- 2  $|A_\epsilon^{(n)}| \leq 2^{n(H(X, Y) + \epsilon)}$
- 3 If  $(\tilde{X}^{(n)}, \tilde{Y}^{(n)}) \sim p(x^{(n)})p(y^{(n)})$ , i.e.,  $\tilde{X}^{(n)}$  and  $\tilde{Y}^{(n)}$  are independent with the same marginals as  $p(x^{(n)}, y^{(n)})$  then

$$P((\tilde{X}^{(n)}, \tilde{Y}^{(n)}) \in A_\epsilon^{(n)}) \leq 2^{-n(I(X; Y) - 3\epsilon)}$$

and for sufficiently large  $n$

$$P((\tilde{X}^{(n)}, \tilde{Y}^{(n)}) \in A_\epsilon^{(n)}) \geq (1 - \epsilon)2^{-n(I(X; Y) + 3\epsilon)}$$

2013-10-09

Joint AEP (see [CT91, Theorem 8.6.1, pp.195-196])

Joint AEP (see [CT91, Theorem 8.6.1, pp.195-196])  
Theorem (Joint AEP)  
Let  $(X^{(n)}, Y^{(n)})$  be sequences of length  $n$  drawn IID according to  $p(x^{(n)}, y^{(n)}) = \prod_i p(x_i, y_i)$ , and choose  $A_\epsilon^{(n)}$  to be the set of jointly typical sequences WRT to  $p(x, y)$  then  
1  $P((X^{(n)}, Y^{(n)}) \in A_\epsilon^{(n)}) \rightarrow 1$  as  $n \rightarrow \infty$   
2  $|A_\epsilon^{(n)}| \leq 2^{n(H(X, Y) + \epsilon)}$   
3 If  $(\tilde{X}^{(n)}, \tilde{Y}^{(n)}) \sim p(x^{(n)})p(y^{(n)})$ , i.e.,  $\tilde{X}^{(n)}$  and  $\tilde{Y}^{(n)}$  are independent with the same marginals as  $p(x^{(n)}, y^{(n)})$  then  
$$P((\tilde{X}^{(n)}, \tilde{Y}^{(n)}) \in A_\epsilon^{(n)}) \leq 2^{-n(I(X; Y) - 3\epsilon)}$$
and for sufficiently large  $n$   
$$P((\tilde{X}^{(n)}, \tilde{Y}^{(n)}) \in A_\epsilon^{(n)}) \geq (1 - \epsilon)2^{-n(I(X; Y) + 3\epsilon)}$$

# Implications of Joint AEP

The jointly typical set has

- about  $2^{nH(X)}$  typical  $X$  sequences
- about  $2^{nH(Y)}$  typical  $Y$  sequences
- about  $2^{nH(X, Y)}$  jointly typical sequences

So when the two variables are not independent,  $H(X, Y) < H(X) + H(Y)$ , and hence not all pairs are jointly typical.

For a fixed  $Y^{(n)}$  we can consider about  $2^{nI(X; Y)}$  such pairs before we are likely to find a jointly typical pair.

That suggests there are about  $2^{nI(X; Y)}$  distinguishable signals  $X^{(n)}$ .

2013-10-09

Implications of Joint AEP

Implications of Joint AEP  
The jointly typical set has  
• about  $2^{nH(X)}$  typical  $X$  sequences  
• about  $2^{nH(Y)}$  typical  $Y$  sequences  
• about  $2^{nH(X, Y)}$  jointly typical sequences  
So when the two variables are not independent,  $H(X, Y) < H(X) + H(Y)$ , and hence not all pairs are jointly typical.  
For a fixed  $Y^{(n)}$  we can consider about  $2^{nI(X; Y)}$  such pairs before we are likely to find a jointly typical pair.  
That suggests there are about  $2^{nI(X; Y)}$  distinguishable signals  $X^{(n)}$ .



# Analysis of Random Codes

## Actual assignment algorithm

- Receiver receives  $Y^n$  according to the distribution

$$P(y^{(n)}|x^{(n)}(w)) = \prod_{i=1}^n p(y_i|x_i(w))$$

- Receiver decodes by guessing that  $w$  is the input that generates a jointly typical sequence:
  - If there is one codeword  $(x^{(n)}(\hat{w}), y^{(n)}) \in A_\epsilon^{(n)}$ , then we decode as  $\hat{w}$ .
  - If there are two codewords such that  $(x^{(n)}(w_i), y^{(n)}) \in A_\epsilon^{(n)}$ , then we declare an error event 2.
  - If there is no codeword  $(x^{(n)}(w), y^{(n)}) \in A_\epsilon^{(n)}$ , then we declare an error event 1.
- In the 1st case, if  $\hat{w} \neq w$  we also declare an error event 2.



Actual assignment algorithm

- Receiver receives  $Y^n$  according to the distribution
 
$$P(y^{(n)}|x^{(n)}(w)) = \prod_{i=1}^n p(y_i|x_i(w))$$
- Receiver decodes by guessing that  $w$  is the input that generates a jointly typical sequence:
  - If there is one codeword  $(x^{(n)}(\hat{w}), y^{(n)}) \in A_\epsilon^{(n)}$ , then we decode as  $\hat{w}$ .
  - If there are two codewords such that  $(x^{(n)}(w_i), y^{(n)}) \in A_\epsilon^{(n)}$ , then we declare an error event 2.
  - If there is no codeword  $(x^{(n)}(w), y^{(n)}) \in A_\epsilon^{(n)}$ , then we declare an error event 1.
- In the 1st case, if  $\hat{w} \neq w$  we also declare an error event 2.

# Analysis of Random Codes

## Probability of errors:

- Probability the jointly typical sequence exists  $\rightarrow 1$  as  $n \rightarrow \infty$  by the first property of the Joint AEP
  - so probability of type 1 errors  $P_1^{error} \rightarrow 0$
- Consider type 2 errors: then for some  $i \neq j$

$$(x^{(n)}(w_i), y^{(n)}(w_j)) \in A_\epsilon^{(n)}$$

- by the code generation process  $x^{(n)}(w_i)$  and  $x^{(n)}(w_j)$  are independent
- hence  $x^{(n)}(w_i)$  and  $y^{(n)}(w_j)$  are independent
- by third property of Joint AEP, for independent  $x^{(n)}(w_i)$  and  $y^{(n)}(w_j)$

$$P((x^{(n)}(w_i), y^{(n)}(w_j)) \in A_\epsilon^{(n)}) \leq 2^{-n(I(X;Y)-3\epsilon)}$$



Probability of errors:

- Probability the jointly typical sequence exists  $\rightarrow 1$  as  $n \rightarrow \infty$  by the first property of the Joint AEP.
  - so probability of type 1 errors  $P_1^{error} \rightarrow 0$
- Consider type 2 errors: then for some  $i \neq j$ 

$$(x^{(n)}(w_i), y^{(n)}(w_j)) \in A_\epsilon^{(n)}$$
  - by the code generation process  $x^{(n)}(w_i)$  and  $x^{(n)}(w_j)$  are independent
  - hence  $x^{(n)}(w_i)$  and  $y^{(n)}(w_j)$  are independent
  - by third property of Joint AEP, for independent  $x^{(n)}(w_i)$  and  $y^{(n)}(w_j)$ 

$$P((x^{(n)}(w_i), y^{(n)}(w_j)) \in A_\epsilon^{(n)}) \leq 2^{-n(I(X;Y)-3\epsilon)}$$

## Analysis of Random Codes

Probability of errors: consider  $w = 1$  WLOG

- There are  $2^{nR}$  codewords, and so  $2^{nR} - 1$  possible incorrect codewords, so the chance of a type 2 error is

$$\begin{aligned} P_2^{error} &\leq (2^{nR} - 1) 2^{-n(I(X;Y) - 3\epsilon)} \\ &\leq 2^{-n(I(X;Y) - 3\epsilon - R)} \end{aligned}$$

- Take rate  $R < I(X; Y) - 3\epsilon$ , then as  $n \rightarrow \infty$  we have

$$P_2^{error} \rightarrow 0$$

This is really just a sketch of a proof, we haven't taken a lot of the nitty little details into account.

Also, we only looked at the 1st part of Shannon's theorem here (that rates up to  $C$  are achievable). Fano's inequality is used to show the converse (that any sequence of codes with error approaching zero must have  $R \leq C$ ) we leave to the reader – see [CT91, Sec8.9].

## Cons of Random Codes

So why don't we use random codes

- 1 very large blocks needed for asymptotic results to hold
- 2 assumes we know  $p(y|x)$
- 3 all codewords must be shared
  - 1  $2^{nR} \times n$  matrix needs to be shared for large  $n$
- 4 decoding very inefficient
  - 1 compute all alternatives and decide which is jointly typical?
  - 2 or store the mapping, which is impractical for even medium blocks

So these random codes are only really suitable for proofs, but there are other places where random codes are used for real, but we will concentrate on some others.

BTW, [CT91] from 1991, says there are no efficient codes that reach capacity – that's not true anymore, just to give an indication of how recent this all is.

## Further reading I

- Thomas M. Cover and Joy A. Thomas, *Elements of information theory*, John Wiley and Sons, 1991.
- David J. MacKay, *Information theory, inference, and learning algorithms*, Cambridge University Press, 2011.