

def.: Let X, Y be R.V. The MUTUAL INFORMATION of X and Y is defined as:

$$I(X, Y) = H(Y) - H(Y|X).$$

$H(Y|X)$ is the conditional entropy of Y given X :

$$H(Y|X) = \sum_{x \in \mathcal{X}} P[X=x] \cdot H(Y|X=x) = \sum_{x \in \mathcal{X}} \left[P[X=x] \sum_{y \in \mathcal{Y}} P[Y=y|X=x] \log \frac{1}{P[Y=y|X=x]} \right].$$

The mutual information is

- symmetric.
- a measure of correlation between X and Y .
- always ≥ 0 .

→ the higher the mutual information, the more X and Y provide information about each other

→ the M.I. between X and Y is minimal ($=0$) when X and Y are independent (since $H(Y|X)=H(Y)$).

We have seen that we cannot compress messages more than their entropy. Nevertheless, in real life, we hardly reach entropy when we communicate. This happens because the communication channel is noisy, and having some redundancy guarantees more ROBUST communication.

→ On one hand, we want to transmit the least number of bits, meaning that we want to compress the information as much as possible (occupying the channel costs, faster communication). At the same time, we want the communication to be reliable, which can be achieved if we add some redundancy in the encoding.

→ The channel allows the communication but at the same time it introduces noise.

Shannon formulated a mathematical model to decode the source of a message while limiting communication errors.

Discrete Memoryless channels (DMC)

High level interpretation: the sender sends a sequence of symbols of a finite alphabet \mathcal{X} , $1 \leq i < \infty$.

The receiver receives a sequence of symbols, possibly from a different finite alphabet \mathcal{Y} , $|\mathcal{Y}| < \infty$. [15]

Since the channel is noisy, some symbols can be altered and the receiver receives a different sequence.

MODEL:

Let $\underline{x} = x_1, \dots, x_n$ be the message sent by the transmitter. In the channel there are $|\mathcal{X}|$ dices, one per symbol. Each dice has $|\mathcal{Y}|$ sides, each with different probabilities.

The probabilities are: $W(\mathcal{Y}|\underline{x}) = P[\text{receive } \mathcal{Y} \text{ | sent } \underline{x}]$.

The probabilities of each symbol in \mathcal{X} and of each transmission are independent of one another, meaning that:

$$W^n(\mathcal{Y}|\underline{x}) = \prod_{i=1}^n W(y_i|x_i) \quad [2]$$

We can define a stochastic matrix $W(\mathcal{Y}|\mathcal{X})$ as:

$$W(\mathcal{Y}|\mathcal{X}) = \begin{pmatrix} W(y_1|x_1) & \dots & W(y_{|\mathcal{Y}|}|x_1) \\ \vdots & & \vdots \\ W(y_1|x_{|\mathcal{X}|}) & \dots & W(y_{|\mathcal{Y}|}|x_{|\mathcal{X}|}) \end{pmatrix} \quad [3]$$

def: A DMC is a channel composed of an input alphabet \mathcal{X} , an output alphabet \mathcal{Y} , $|\mathcal{X}|, |\mathcal{Y}| < \infty$ and a transition stochastic matrix $W(\mathcal{Y}|\mathcal{X})$ as the one in [3] that satisfies [2].

def: An ENCODER $\mathcal{C}^n \subseteq \mathcal{X}^n$ is a set of distinct sequences of length n that the transmitter can communicate through a channel.

- it maps messages to "codewords". $|\mathcal{C}^n| = \# \text{ of different messages that can be transmitted through the channel}$.

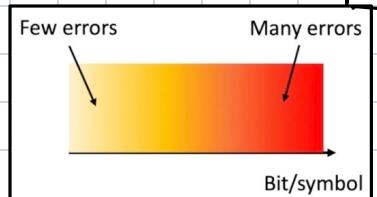
def: A DECODER $q_n : \mathcal{Y}^n \rightarrow \mathcal{C}^n$ is a function that associates received messages of length n to sequences of \mathcal{C}^n .

To be good, the communication must be s.t. the probability that a received message \mathcal{Y} is associated to the correct message \underline{x} is high. At the same time, we want to compress the information as much as possible, so that we can send more information

↳ We want to maximise $\sum_{\mathcal{Y} \in q_n^{-1}(\underline{x})} W^n(\mathcal{Y}|\underline{x})$ (= probability of receiving the correct messages) and $|\mathcal{C}^n|$.

The problem is that these two quantities have inverse trends!

If each sent symbol carries a lot of information (expressed in bits), meaning that the information is strongly compressed, then an error occurring on the channel can cause the misunderstanding of the message.



Shannon's theorem, stated later, establishes the relationship between how much data can be sent on a MDC with an almost-zero error probability.

- ! GOAL = maximise the speed of the transmission (i.e. the size of C^n that can be transmitted in a unit of time), while keeping very low errors.
- HOW CAN WE DO IT? We CANNOT change W , that is an intrinsic characteristics of the channel. The only element that we can tune is C^n through $\underline{P(x)}$

def.: The Transmission rate of an encoder C^n is the number of bits per transmission

$$\frac{1}{n} \lg_2 |C^n| \rightarrow \text{it's the number of bits carried by a symbol.}$$

An error occurs when the transmitter sends a sequence x , and the receiver receives a sequence y s.t. $\varphi(y) \neq x$, or, equivalently, if $y \notin J_m^{-1}(x)$

example. $X = \{0, 1\}$, $Y = \{a, b\}$. $\varphi_1 : Y \rightarrow X$: $\varphi_1(a) = 0, \varphi_1(b) = 1$

$$W(y=a|x) \quad W(y=b|x)$$

$$W(y|x) = \begin{cases} W(y|x=0) & \boxed{0.8} \\ W(y|x=1) & \boxed{0.2} \end{cases} \quad \begin{matrix} \text{CORRECT} \\ \text{WRONG} \end{matrix}$$

$$X = (x_1, \dots, x_n), x_i \in \mathcal{X}, \quad Y = (y_1, \dots, y_n), y_i \in \mathcal{Y}.$$

$$\varphi_n : Y^n \rightarrow X \text{ s.t. } \varphi_n(y) = (\varphi_1(y_1), \dots, \varphi_n(y_n))$$

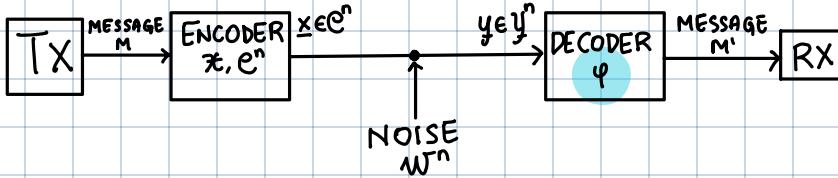
$$\text{If } n=2, \text{ then, } C^2 = \{00, 01, 10, 11\},$$

$\varphi_2(00) = aa, \varphi_2(01) = ab, \varphi_2(10) = ba, \varphi_2(11) = bb.$

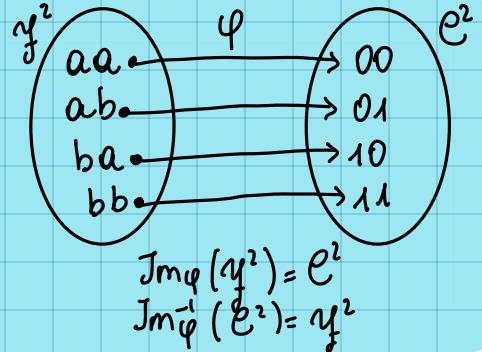
$P[aa|00] = W^2(0|0) = 0.8^2. \quad P[ab|00] = 0.8 \cdot 0.2$

IF the sender sends $\underline{x} = 00$ and the receiver receives $\underline{y} = aa$, no errors occurred, because $aa \in Jm_{\varphi}^{-1}(00)$ [17]

IF the sender sends $\underline{x} = 00$ and the receiver receives $\underline{y} = ab$, some error occurred, since $ab \notin Jm_{\varphi}^{-1}(00)$



$$\begin{aligned} Jm_{\varphi}(\{aa\}) &= \{00\}, Jm_{\varphi}^{-1}(\{00\}) = \{aa\} \\ Jm_{\varphi}(\{ab\}) &= \{01\}, Jm_{\varphi}^{-1}(\{01\}) = \{ab\} \\ Jm_{\varphi}(\{ba\}) &= \{10\}, Jm_{\varphi}^{-1}(\{10\}) = \{ba\} \\ Jm_{\varphi}(\{bb\}) &= \{11\}, Jm_{\varphi}^{-1}(\{11\}) = \{bb\} \end{aligned}$$



def: the error probability for e^n and φ^n is:

$$W^n(\overline{\varphi_n^{-1}} | \underline{x}) = 1 - W^n(\varphi_n(\underline{x}) | \underline{x})$$

$\underbrace{\qquad}_{\downarrow}$

this means $y^n - Jm_{\varphi}^{-1}(\underline{x})$ [difference between sets]

In the previous example, if $\underline{x} = 00$, $y^2 - Jm_{\varphi}^{-1}(00) = \{ab, ba, bb\}$.

$$W^n(\overline{\varphi_2^{-1}} | \underline{x}) = 1 - W(a|0)^2 = 1 - 0.8^2 = 0.36$$

def: the maximum error probability is:

$$e_n(W^n, C^n, \varphi^n) = \max_{x \in e^n} W^n(\overline{\varphi_n^{-1}(x)} | \underline{x})$$

def: REN is an ACHIEVABLE TRANSMISSION RATE for the DMC (W^n, x, y) if \exists sequence $\{C^n, \varphi^n\}_{n \in \mathbb{N}}$ s.t.

$$e_n(W^n, C^n, \varphi^n) \rightarrow 0 \quad \lim_{n \rightarrow \infty} \frac{1}{n} \lg |C^n| \geq R$$

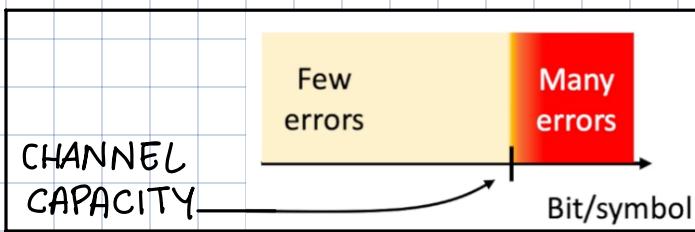
An achievable transmission rate is the number of bits that can be transmitted in a channel with an error that goes to zero.

OBS: all achievable rates are FINITE, $R < \infty$.
This means that \exists maximum achievable rate.

def: the capacity of a DMC W is the maximum achievable rate:

$$C(W) = \max_R \{ R : R \text{ is an achievable rate} \}$$

So, basically, the capacity of a channel is the MAXIMUM ACHIEVABLE RATE (NUMBER OF BITS) THAT CAN BE TRANSMITTED WITH AN ERROR PROBABILITY THAT GOES TO ZERO.
IT IS THE MAXIMUM RELIABLE TRANSMISSION RATE.



def: The Shannon's capacity of a channel W is:

$$C^*(W) = \max_{(X,Y) \sim W(Y|X)} I(X, Y)$$

that is, the maximal mutual information between sent messages and received messages.

SHANNON'S CAPACITY THEOREM: $C(W) = C^*(W)$

This result is very powerful. It says that the capacity of a channel depends on how much information can be transmitted reliably over a noisy channel.

The model that we have seen decouples the communication process into two tasks:

- SOURCE CODING (data compression) \rightarrow Reduce # of bits to represent data
- CHANNEL CODING (error correction) \rightarrow Add redundancy to make transmission reliable.

Continuous channels

19

We've spoken about DISCRETE CHANNELS.

In IoT, we usually deal with CONTINUOUS, WIRELESS CHANNELS

→ The Shannon capacity theorem has an alternative formulation in this case.

We consider a channel with  Gaussian noise and a given bandwidth, B . The noise is independent of the signal.

- Gaussian noise.

→ We will see that channels are subject to ADDITIVE INTERFERENCES (SPOILER: WIRELESS communication happens through electromagnetic waves) that result in adding noise to the signal. Such noise usually follows a Gaussian or NORMAL distribution, which has probability density function

$$q(z) = \frac{e^{-\frac{(z-\mu)^2}{2\sigma^2}}}{\sqrt{2\pi}}$$

where z is the noise value, μ is the mean and σ is the standard deviation

If x is the input signal and y is the output signal, then we have:

$$y = x + z, \quad z \sim N(0, \sigma_z)$$

→ We need to find $C^*(w) = \max_{(Y,X)} I(X,Y)$
 $P_{Y|X} = w$

$$I(X,Y) = H(Y) - H(Y|X).$$

$$H(Y|X) = H(X+z|X) \text{ but since } X \perp\!\!\! \perp z, H(X+z|X) = H(z)$$

$$\text{Hence } I(X,Y) = H(Y) - H(z).$$

What is the $P_{Y|X}$ s.t. $I(X,Y)$ is maximal?

→ When also $x \sim N(0, \sigma_x)$! [proof omitted].

There is a theorem in probability that says that if x, z are normally distributed R.V. and are independent, and have variance σ_x^2 and σ_z^2 respectively, then: $x+z$ is also normally distributed and its variance is $\sigma_{x+z}^2 = \sigma_x^2 + \sigma_z^2$

The problem is that we have always spoken about discrete random variables, and we have defined entropy for discrete probability distributions.

Good news is that we can define entropy for continuous R.V.s too! [differential entropy]

Without digging into details, the entropy for a normally-distributed R.V. with variance σ is

$$\frac{1}{2} \lg (2\pi e \sigma^2)$$

Euler constant

$$\begin{aligned} \text{So: } C^*(W) &= \frac{1}{2} \lg (2\pi e \sigma_y^2) - \frac{1}{2} \lg (2\pi e \sigma_z^2) = \\ &= \frac{1}{2} \lg \left(\frac{2\pi e \sigma_y^2}{2\pi e \sigma_z^2} \right) = \frac{1}{2} \lg \left(\frac{\sigma_y^2}{\sigma_z^2} \right) = \frac{1}{2} \lg \left(\frac{\sigma_x^2 + \sigma_z^2}{\sigma_z^2} \right) = \\ &= \frac{1}{2} \lg \left(1 + \frac{\sigma_x^2}{\sigma_z^2} \right) = \begin{cases} \text{The VARIANCE of a signal is equal} \\ \text{to its POWER. THE QUANTITY} \\ \frac{\sigma_x^2}{\sigma_z^2} \text{ is the SIGNAL TO NOISE RATIO} \\ (\text{SNR}) \end{cases} \\ &= \frac{1}{2} \lg (1 + \text{SNR}) \end{aligned}$$

To far, we have considered only the transmission of one signal x . But we can send a sequence of signals. How long the sequence of signals can be is determined by the bandwidth of the channel, B . Hence we get the

SHANNON-HARTLEY THEOREM $C^*(W) = B \lg (1 + \text{SNR})$

(Follows from the NYQUIST-SHANNON theorem. We have mentioned it in our 4th lecture - EMBEDDED SYSTEMS).