

1 INTRODUCTION:

Nowadays, manipulating daily loads of info: images, text, sounds = files = 0.11000 Gb which we are performing: - compression $\text{file} \rightarrow \text{smaller file}$

- transmission over noisy channel $\text{file} \xrightarrow{\text{noise}} \text{computer file}$
error correction

These two problems were already formalised and approximately solved by Shannon 48 "A mathematical theory of communication".

Why should we care about this while not being telecommunication engineer?

bio applicat.: - DNA \rightarrow RNA \rightarrow proteins

- vision: receptors/ retina \rightarrow optical nerve \rightarrow brain

social: - language eg: - repetition \rightarrow error correcting.

- only some combination of sounds makes up words

theoretically: Information theory enlightens statistical mechanics.

\rightarrow faithful interdisciplinary of stat mech of disordered systems. (Cris, Florent, Federico)
eg. Constraint satisfaction problems for error correcting LDPC

\rightarrow WEB PAGE: www.phy.ens.fr/~guilhem/boulder.html

2 THE MEANING OF ENTROPY:

2A definition: $S_{mc} = k_B \ln \Omega$
microcanonical

$S_C = -k_B \sum_c p(c) \ln p(c)$ with $p(c) = \frac{1}{Z} e^{-\beta H(c)}$

More I.T:

\hookrightarrow Shannon entropy for distribution $p = \{p(x); x \in X\} \equiv H(p) = -\sum_{x \in X} p(x) \log_2 p(x)$
(change of unit of Boltzmann k_B). \downarrow alphabet

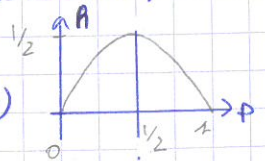
prop: $H(p) = 0 \Leftrightarrow p(x) = \delta_{x, x_0}$ concentrated on one value
maximal

$H(p) \in [0, \log_2 |X|]$ and $H(p) = \log_2 |X| \Leftrightarrow p(x) = \frac{1}{|X|} \forall x$

"the entropy measures the lack of information on a realization of p "

\hookrightarrow yet there are many functions growing between uniform and concentrated, so why should it be $H(p)$ with the logarithm?

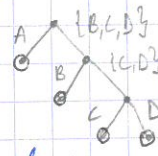
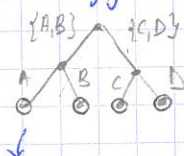
definition for bernoulli(p): $h(p) = -p \log_2 p - (1-p) \log_2 (1-p)$
 $p \in (0, 1]$



2B guessing game: } Amy choose $x \in X$

} Sheldon has to guess x asking YES/NO questions as fast as possible.
a strategy: another one:


eg. $X = \{A, B, C, D\} \rightarrow$



formalisation: $T = \text{strategy} = \text{tree of questions}$
 $l_x(T) = \# \text{ questions to find } x \in X \text{ with strategy } T$

to define the best strategy we compute $\bar{E}(T) = \sum_{x \in X} p(x) l_x(T) \Rightarrow T^* = \underset{T}{\text{argmin}} \bar{E}(T)$

this is bounded: $H(p) \leq \bar{E}(T^*) \leq H(p) + 1$ SHANON SOURCE CODING TRFT
 measures the lack of information within one bit of x knowing p !!
 " T^* minimises the number of questions Sheldon expects to ask"

a proof: - Kraft inequality $\forall T, \sum_{x \in X} 2^{-l_x(T)} \leq 1$
 $l_{\max} = \max_{x \in X} l_x(T)$
 total number leaves at l_{\max} : $2^{l_{\max}}$
 total number of shadows: $\sum_{x \in X} 2^{l_{\max} - l_x}$

 $\rightarrow \text{Kraft.}$

- Kullback-Leibler divergence: p, q proba laws on X $D(p||q) = \sum_{x \in X} p(x) \log_2 \frac{p(x)}{q(x)}$
 $p(x) \in (0, 1]$
 $\sum p(x) = 1$
 prop: $D(p||q) \geq 0$ (due to the concavity of the log)
 $D(p||q) = 0 \Leftrightarrow p = q$ } CAN BE THOUGHT AS A DISTANCE

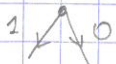
- the proof:
 * consider $q(x) = \frac{1}{Z}$ with $Z = \sum_{x \in X} 2^{-l_x} \leq 1$ by Kraft
 $\hookrightarrow D(p||q) = \sum_x p(x) \log_2 \frac{p(x)}{2^{-l_x}} Z = -H(p) + \log_2 Z + \bar{E}(T) \geq 0$
 $\Rightarrow H(p) \leq \bar{E}(T)$

* if $\{l_x\}_{x \in X}$ integers s.t. $\sum_x 2^{-l_x} \leq 1 \Rightarrow \exists T$ such that $l_x(T) = l_x \forall x$
 $l_{x_1} \leq l_{x_2} \leq \dots$ put x_1 in the leftmost node of depth l_{x_1}
 x_2 in the remaining leftmost node of depth l_{x_2}

we would like to take $l_x = \lceil \log_2 \frac{1}{p(x)} \rceil \rightarrow$ obey Kraft inequality
 $\leq \log_2 \frac{1}{p(x)} + 1$
 $\bar{E}(T) \leq \sum_x p(x) \lceil -\log_2 p(x) + 1 \rceil \leq H(p) + 1$ ET!

2C data compression:

consider a file of string of symbols in X : $x_1 x_2 x_3 \dots \rightarrow 00101\dots$ as short as possible

It is actually the problem we have solve:  so that $A=0$
 $B=10$
 $C=110$
 $D=111$

let's formalise that with $w: x_i \rightarrow w_{x_i}$
 stacked strings of words built in these way are uniquely decodable, i prefix free
 going back to the root until reaching a word (over and over) - m. the fl. 1

if x_i 's are iid with proba $p(x)$ → need $nH(p)$ bits to compress x_1, \dots, x_n
 \equiv "entropy is the best rate of compression possible".

2D mutual information

An inference problem: given n obs of random X from observation Y .

formally: the pair of r.v $(X, Y) \in (\mathcal{X}, \mathcal{Y})$

distributed according to $p_{X,Y}(x,y) = P[X=x, Y=y]$

→ marginal laws: $p_X(x) = \sum_y p_{X,Y}(x,y)$, $p_Y(y) = \sum_x p_{X,Y}(x,y)$

→ conditional laws: $p_{X|Y}(x|y) = \frac{p_{X,Y}(x,y)}{p_Y(y)}$

From which we consider various interesting entropies:

* joint law entropy $H(X,Y) = - \sum_{x,y} p_{X,Y}(x,y) \log_2 p_{X,Y}(x,y)$

* entropy of marginal $H(X) = - \sum_x p_X(x) \log_2(p_X(x))$

* conditional ent. $H(X|Y) = - \sum_y p_Y(y) \sum_x p_{X|Y}(x|y) \log_2 p_{X|Y}(x|y)$
 $= - \sum_{x,y} p_{X,Y}(x,y) \log_2 p_{X|Y}(x|y)$

* mutual info $I(X,Y) = D(p_{X,Y} \parallel p_X p_Y) \Rightarrow$ prop: $I(X,Y) = 0 \Leftrightarrow X, Y$ indep.

$$= \sum_{x,y} p_{X,Y}(x,y) \log \frac{p_{X,Y}(x,y)}{p_X(x)p_Y(y)}$$

$$= H(X) + H(Y) - H(X,Y) \Rightarrow$$
 prop: $H(X,Y) \leq H(X) + H(Y)$

$$= \sum_{x,y} p_{X,Y}(x,y) \log \frac{p_{X,Y}(x,y)}{p_X(x)}$$

$$= H(X) - H(X|Y) \Rightarrow$$
 prop: $H(X|Y) \leq H(X)$

"conditioning reduce entropy"

- EXERCISES: - show that Gibbs maximize entropy under constraint $\langle E \rangle = U$
 - prove that uniform distribution maximizes entropy under no-constraint (KL div)
 - SHANNON'S PAPER → Appendix d:

$$\begin{cases} H(p_1, \dots, p_M) \text{ continuous in } p_i \\ H(\frac{1}{n}, \dots, \frac{1}{n}) \text{ growing with } n \end{cases}$$

$$\Rightarrow H = - \text{const} \sum_{i=1}^n p_i \log p_i$$

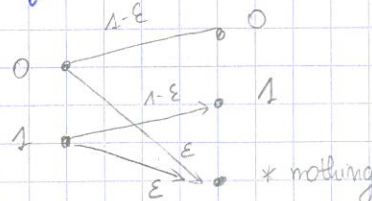
$$H\left(\begin{matrix} p_1 \\ p_2 \\ p_3 \end{matrix}\right) = H\left(\begin{matrix} p_1 \\ p_2+p_2 \end{matrix}\right) + H\left(\begin{matrix} p_2/p_2+p_3 \\ p_3/p_2+p_3 \end{matrix}\right)$$

3 COMMUNICATION OVER NOISY CHANNELS

3A definitions message $\xrightarrow{\text{noise}}$ corrupted message.

examples: → Binary Erasure Channel: (BEC)

$\varepsilon \equiv$ probability of erasure.



→ Binary Symmetric Channel: (BSC)



in which case you can never be entirely sure

10/07/2017

"the higher, the better you can deduce input from output"

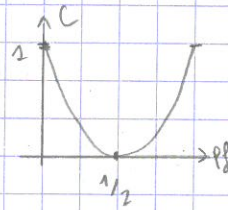
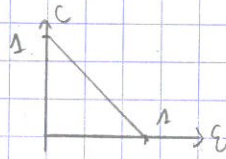
Capacity of a channel:

$$C = \max_{P_X} I(X; Y) \text{ with } X \text{ input of a channel.}$$

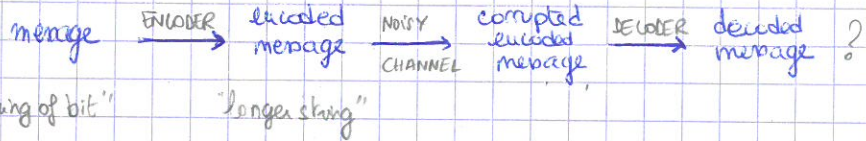
Y output
 P_X probability law of input

EXERCISES: $C_{BEC} = 1 - \epsilon$

$$C_{BSC} = 1 - h(p_f)$$



Encoding and decoding:



Rate of a code: $\equiv \frac{\text{\# bits of message}}{\text{\# bits of encoded msg.}} < 1$ → the larger the better to reduce the cost of message transmission.

3B naive coding \equiv repetition

e.g. we repeat 3 times input:

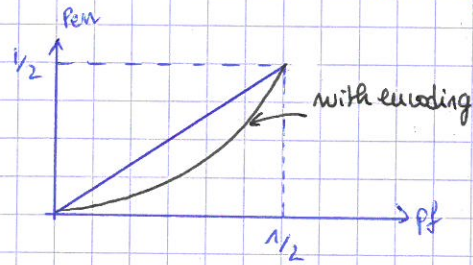
ENCODER	0	→	0 0 0	with the BSC $p_f < 1/2$.
	1	→	1 1 1	

reasonable DECODER "majority rule" → add number of bits

000	→	0	rate = 1/3
111	→	1	
001, 010, 100	→	0	
110, 101, 011	→	1	

How good is the naive coding?

Probability of error without encoding = p_f
 Probability of error with 3 repetitions = $p_f^3 + 3p_f^2(1-p_f)$
3 flips $\binom{3}{2}$ flips



→ better than no encoding but $P_{en} > 0$ as soon as $p_f > 0$
 rate = 1/3
 → actually not that good. Can one do better?

3C Shannon channel coding theorem

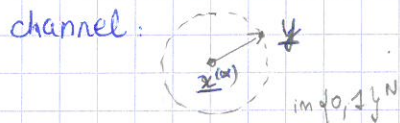
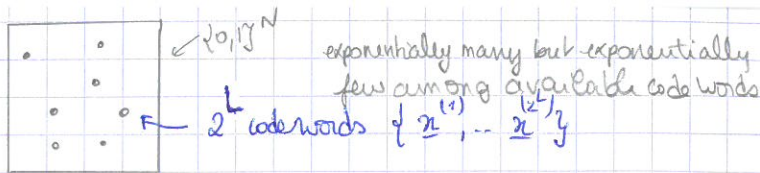
THM: There exist encoding (of growing length) with $P_{en} \xrightarrow{N \rightarrow \infty} 0$, for all rates R smaller than the capacity of the channel C ($R < C$).

So that we can now interpret the capacity as the best achievable rate with $P_{en} \rightarrow 0$.
 The paradoxical statement that we could be sure of the signal despite the noise is resolved by the fact this is an asymptotic statement (thermodynamic limit).

MORE FORMAL DEFINITIONS:

message $\underline{x} = (x_1, \dots, x_N)$ $\in \{0, 1\}^N$	encoded message $\underline{z} = (z_1, \dots, z_N)$ $\in \{0, 1\}^N$	corrupted eval. msg. $\underline{y} \in X_{out}^N$ BEC $X_{out} = \{0, 1, N\}$ BSC $X_{out} = \{0, 1\}$	corrected encoded msg. $\hat{\underline{z}}(\underline{y})$	corrected msg. $\hat{\underline{x}}(\hat{\underline{z}}(\underline{y}))$
--	--	--	--	---

encoding: $\underline{x} = f_{\text{encoding}}(\underline{z})$
 Code book = $\mathcal{C} = \{ \underline{x}^{(1)}, \dots, \underline{x}^{(2^L)} \}$



intuition \rightarrow put the codewords as far as possible in space
 BUT! CAREFUL IN HIGH DIMENSIONS, things do not happen the same way as in 2d!! oversimplifying view-

correction: $\hat{\underline{x}}(\underline{y})$ estimation of $\underline{x}^{(a)}$
 decoding: $\hat{\underline{z}} = f_{\text{decoding}}^{-1}(\hat{\underline{x}}(\underline{y}))$

\rightarrow From now on we will focus on the central $\underline{x} \rightarrow \underline{y} \rightarrow \hat{\underline{x}}(\underline{y})$. If the $2^L \underline{z}$ are uniform at random, the \underline{x} are also u.a.r.

DECODING AS AN INFERENCE PROBLEM:

\underline{x} random code word, $P_{\underline{x}}(\underline{x}) = \frac{1}{2^L} \mathbb{1}(\underline{x} \in \mathcal{C})$
 \underline{y} random output $P_{\underline{y}|\underline{x}}(\underline{y}|\underline{x}) = \prod_{i=1}^N Q(y_i|x_i)$
 \uparrow over the bits indep.

BEC: $\begin{cases} Q(0|0) = 1 - \epsilon \\ Q(1|0) = \epsilon \\ Q(1|1) = 0 \end{cases}$

using Bayes theorem: $P_{\underline{x}|\underline{y}}(\underline{x}|\underline{y}) = P_{\underline{y}|\underline{x}}(\underline{y}|\underline{x}) \frac{P_{\underline{x}}(\underline{x})}{P_{\underline{y}}(\underline{y})}$ \leftarrow prior proba on the signal

$$= \prod_i Q(y_i|x_i) \frac{1}{2^L} \mathbb{1}(\underline{x} \in \mathcal{C}) \frac{1}{P_{\underline{y}}(\underline{y})}$$

$$= \frac{1}{Z(\underline{y})} \mathbb{1}(\underline{x} \in \mathcal{C}) \prod_{i=1}^N Q(y_i|x_i) \quad \text{POSTERIOR PROBABILITY}$$

One keeps only the factors function of \underline{x} and lets the rest into the partition $Z(\underline{y})$.

RR: Having all the distribution of the decoding given the observation implies that we have \neq ways of decoding:

- * $\hat{\underline{z}}(\underline{y}) = \underset{\underline{z}}{\text{argmax}} P_{\underline{x}|\underline{y}}(\underline{z}|\underline{y})$ block MAP $\rightarrow \min P(\hat{\underline{x}} \neq \underline{x})$ (Maximal a posteriori -)
- * $\hat{\underline{x}}_i(\underline{y}) = \underset{x_i}{\text{argmax}} P_{x_i|\underline{y}}(x_i|\underline{y})$ symbol MAP decoding $\rightarrow \min E(d(\underline{\hat{x}}, \underline{x}))$

EXAMPLE WITH THE BEC: $\underline{y} = (0, 1, 0, 0, *, *, 0, \dots, *, 0)$

$$P_{\underline{x}|\underline{y}}(\underline{x}|\underline{y}) = \frac{1}{|\mathcal{E}(\underline{y})|} \mathbb{1}(\mathcal{E}(\underline{y}) \cap \mathcal{C}) \quad \mathcal{E}(\underline{y}) = \{ \underline{x} \in \{0,1\}^N \text{ s.t. } y_i \in \{0,1\} \Rightarrow x_i = y_i \}$$

\leftarrow intersection of codebook and ball

uniform probability over the messages matching the revealed bits in received message -

How should we construct code books? Let's start with a simple construction proposed by Shannon.

SHANNON RANDOM CODE ENSEMBLE \rightarrow connection to the RRC!

$\mathcal{C} = \{ \underline{x}^{(1)}, \dots, \underline{x}^{(2^L)} \}$ with $\underline{x}^{(a)} = \{ x_1^{(a)}, \dots, x_N^{(a)} \} \rightarrow$ choose the $N \cdot 2^L$ $x_i^{(a)}$ and 0,1 with probability $1/2$.

RR: proba that $\underline{x}^{(A)} = \underline{x}^{(B)}$ for $A \neq B$ (non-injective)

$\xrightarrow[N \rightarrow \infty]{L \rightarrow \infty} 0$
 $R = \frac{L}{N}$ fixed

Bin(N, 1-ε) fluctuations

ANALYSIS ON THE BEC:

Assume w.p. 1-ε that $\underline{x}^{(1)}$ has been sent: $\underline{y} = \underline{x}^{(1)}$ on $(N(1-\epsilon) + \mathcal{O}(1))$ correctly transmitted
 \times on other bits $\sim \text{NE}$

$W = \sum_{i=1}^N \text{Bin}(2^L - 1, (\frac{1}{2})^{N(1-\epsilon)})$ "number of confusing code word"

$\Rightarrow E(W) = \sum_{i=1}^N \frac{2^L - 1}{N(1-\epsilon)} \rightarrow 0$ if $\epsilon < 1 - R$ $\Rightarrow P_r[W \neq 0] \rightarrow 0$ -
total # of available code words probab 2 + codewords agree on a bit
integers!

First moment / Markov inequality for strictly positive code word -

So that no confusion if $R < 1 - \epsilon = C_{BEC} \rightarrow$ Shannon theorem!

Even the random codes achieve capacity - (here we considered averages for codewords etc...)

HOMEWORK: proof for the BSC (text on website) -

Rk we prove that there exist such codes, but we did not show that we could not do better than this bound - if $R > C \Rightarrow P_{err} > 0$ coming from $H(\underline{X}|\underline{Y}) > 0$, Fano inequality -

⚠ In practice, encoding and decoding the Random Code ensemble takes exponential time and memory - \mathcal{C} has to be described by $N \cdot 2^L$ bits. The encoding is easy by looking up the table. But the decoding is NP hard as one need to look for the corresponding match in the table - The randomness hinders any compression, we can only do exhaustive search in the table -

4 LOW DENSITY PARITY CHECK CODES (LDPC). \rightarrow putting some structure

4A Linear codes

- $\{0, 1\}^N$ is a linear space over $\mathbb{Z}_2 = \{0, 1\}$ "scalars", addition mod 2, multiplication
- The codebook $\mathcal{C} \subset \{0, 1\}^N$ is said to be a linear code if \mathcal{C} is a linear subspace of $\{0, 1\}^N$.
 $\underline{x} + \underline{y} = (x_1 + y_1, x_2 + y_2, \dots, x_N + y_N) \in \mathcal{C}$
mod 2
- $\underline{x} + \underline{x} = \underline{0} \in \mathcal{C}$ (the origin belongs to the linear subspace).
- $\hookrightarrow \underline{0}$ is always a codeword of a linear code
 {All codewords are equivalent (can always make gauge transformation to change position of origin)}.