

# Bridges between statistical mechanics and information theory

Guilhem  
Semejian

## I Introduction

why these lectures, why in this school

- information everywhere, "big data" era

↳ texts, sounds, images... = computer files = 0110...

2 basic problems: . compressing these files ← how much information in it? (no moral prejudice)  
. transmitting through noisy channels (interference of cell phones)

formalized by Shannon, 48, A mathematical theory of communication, fundamental paper

- other forms of information treatment

- DNA → RNA → <sup>translation</sup> proteins → cells → organs → organisms, lot of noise, regulations  
<sub>transcription</sub>

- photons → retina → optical nerve → brain perception

- language, possibility of error corrections

- relationship with the school

• enlightens fundamental features of stat mech (entropy)

• one of the fruitful interdisciplinary applications of stat mech of disordered systems (noise = disorder)  
↳ PC, CSPs

→ relations with  
C. Moore  
G. Biroli  
H. Cohen (packing in Hamming space)  
F. Kuzakala  
F. Ricci-Tersaghi

deviations from original title of the lecture, but we'll see at the end  
a few of "cavity computation", to be completed by Federico

bibliography,

outline of the lectures,

exercises on:

[www.phys.ens.fr/~guilhem/boulder.html](http://www.phys.ens.fr/~guilhem/boulder.html)

## II The meaning of entropy

### II. A. Definition

in stat. mech.  $S = k \ln \Omega$  microcanonical

$$S = -k \sum_{\mathcal{E}} p(\mathcal{E}) \ln p(\mathcal{E}) \quad \text{canonical, } p(\mathcal{E}) = \frac{e^{-\beta H(\mathcal{E})}}{Z}$$

if  $p(\mathcal{E}) = \begin{cases} \frac{1}{\Omega} & \text{on } \Omega \text{ configurations} \\ 0 & \text{o.w.} \end{cases}$

reduces to microcanonical entropy

Shannon's definition and notation:

$p = \{ p(x), x \in X \}$  probability law

$$H(p) = - \sum_{x \in X} p(x) \log_2 p(x)$$

choice of units of  $k = \frac{1}{\ln 2}$

$$0 \ln 0 = 0$$

"entropy is a measure of (the lack of) information"

$$p(x) = \delta_{x, x_0} \iff H(p) = 0$$

$$p(x) = \frac{1}{|X|} \iff H(p) = \log_2 |X| \quad \text{maximal (exercise, to be done with } \mathcal{D}(p||q) \text{)}$$

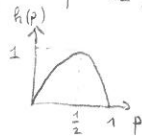
and in general  $H(p) \in [0, \log_2 |X|]$

$\Rightarrow$  grows when  $p$  spreads out, ie when randomness  $\uparrow$ , but

why this precise form? many other functions could be used

we shall see that this is the "right" def.

if  $|X|=2$ , 
$$h(p) = -p \log_2 p - (1-p) \log_2 (1-p)$$

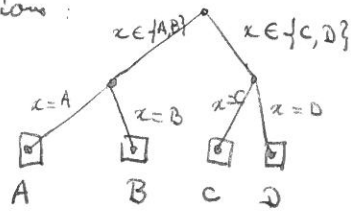


## II.B. A guessing game

2 players, Alice draw r.v.  $X \in \mathcal{X}$  (proba  $p(x) = \mathbb{P}(X=x)$ )  $\rightarrow$  outcome  $x$   
to be precise later  
 Bob does not see  $x$ , ask yes/no questions to Alice to determine  $x$   
Sheldon

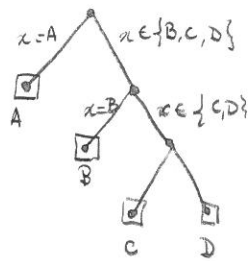
example:  $\mathcal{X} = \{A, B, C, D\}$

tree of questions:



$l_A = l_B = l_C = l_D$ ,  $l_x$ : nb of questions to determine  $x$

but Bob can also play with:



$l_A = 1$ ,  $l_B = 2$ ,  $l_C = l_D = 3$

what's the best strategy? ie for Bob to conclude as fast as possible  
 it depends on  $p(x)$ , if  $A$  is very probable,  $2^1$  is better

call  $T$  the tree of questions,  $l_x(T)$  the nb of questions to conclude  $x$  in  $T$

$$\bar{l}(T) = \sum_x p(x) l_x(T) \quad \text{average nb of questions asked}$$

$T^*$  the choice that minimizes  $\bar{l}(T)$

claim (Shannon's source coding th):

$$H(p) \leq \bar{l}(T^*) < H(p) + 1$$

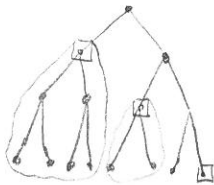
$\Rightarrow H(p)$  really quantifies the lack of information on the output of  $p$ ,  
 as the minimal nb of binary questions to ask to determine the outcome

elements of the proof:

• Kraft inequality:  $\forall T, \sum_{x \in X} 2^{-l_x(T)} \leq 1$

ie if  $|X|$  grows, some  $l_x(T)$  must be big to compensate

proof:  $l_{max} = \max_x l_x$



$2^{l_{max}}$  nodes at depth  $l_{max}$

each  $l_x$  projects a "shadow" on  $2^{l_{max} - l_x}$  nodes

shadows do not intersect  $\Rightarrow \sum_{x \in X} 2^{l_{max} - l_x} \leq 2^{l_{max}}$

by definition  $\square \Rightarrow$  answer found  $\Rightarrow$  no other  $\square$  below  
 divide by  $2^{l_{max}} \Rightarrow$  done

• Kullback Leibler divergence

$p, q$  two proba. laws on  $X$ , def  $D(p||q) = \sum_{x \in X} p(x) \log_2 \frac{p(x)}{q(x)}$

properties:  $D(p||q) \geq 0$  :  $D = - \sum_{x \in X} p(x) \log_2 \frac{q(x)}{p(x)}$   $\log_2$  concave

$\leq \log_2 \left( \sum_{x \in X} p(x) \frac{q(x)}{p(x)} \right) = \log_2(1) = 0$   
Jensen

\*  $D(p||q) = 0 \Rightarrow p(x) = q(x) \forall x \in X$

because  $\log$  strictly concave, Jensen saturated  $\Rightarrow f^o = \text{at}$

rk: with  $q(x) = \frac{1}{|X|}$   $D(p||q) = -H(p) + \log_2(|X|) \Rightarrow H(p) \leq \log_2 |X|$ , equality only if  $p=q$

proof of Jensen:  $E[f(x)] \leq f(E[x])$  for  $f$  concave

$X = E[X] + (X - E[X])$

$f(a+b) \leq f(a) + (b-a)\lambda$

$f(x) \leq f(E[X]) + (x - E[X])\lambda$

$E[f(x)] \leq f(E[X]) + 0 \cdot \lambda$

$\lambda = f'(a)$  if  $f$  derivable,  
 even if not derivable such a  $\lambda$  exists

• putting together the two

$$q(x) = \frac{1}{z} 2^{-l_x(T)} \quad \text{for a valid tree,} \quad z = \sum_{x \in X} 2^{-l_x(T)}$$

$$D(p||q) = \sum_x p(x) \log_2 \left( \frac{p(x)}{\frac{1}{z} 2^{-l_x(T)}} \right) = -H(p) + \log_2(z) + \bar{l}(T)$$

$$\text{Kraft} \Rightarrow z \leq 1 \quad \log_2(z) \leq 0$$

$$\bar{l}(T) - H(p) = \underbrace{D(p||q)}_{\geq 0} - \underbrace{\log_2(z)}_{\geq 0} \geq 0 \quad \Rightarrow \quad H(p) \leq \bar{l}(T) \quad \forall \text{ valid } T$$

$$\Rightarrow H(p) \leq \bar{l}(T^*)$$

this proves the lower bound of Shannon's th

upper bound: one can exhibit a tree which achieves it:

\* if  $\{l_x\}$  set of integers satisfying Kraft, then  $\exists T$  with  $l_x(T) = l_x$ :  
 order the  $l_x$   $l_1 \geq l_2 \geq \dots$ , take the first node in lexicographic order of depth  $l_1$ , then remove what is below and continue

$$* \text{ set } l_x = \left\lceil \log_2 \frac{1}{p(x)} \right\rceil \geq -\log_2 p(x), \quad \sum_x 2^{-l_x} \leq \sum_x p(x) = 1 \Rightarrow \text{Kraft,}$$

$$\sum_x p(x) l_x \leq \sum_x p(x) \left( \log_2 \frac{1}{p(x)} + 1 \right) = H(p) + 1$$

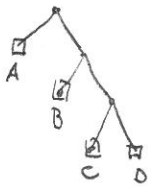
Huffman coding to find the optimal choice of  $\{l_x\}$  given  $p(x)$ : group the two least likely symbols to make a single one, and iterate

II. C. Data compression

back to original problem :  $x_1 x_2 \dots x_m \rightarrow 01001 \dots 01$   
 ↑  
 EX  
 as short as possible  
 but that allows to recover

we have actually solved this problem

tree



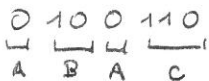
left branch  $\leftrightarrow 0$   
 right branch  $\leftrightarrow 1$

- A: 0
- B: 10
- C: 110
- D: 111

each symbol  $x$  associated to a string  $w_x(T)$  of 0,1, of length  $l_x(T)$

$$x_1 \dots x_m \rightarrow w_{x_1}(T) w_{x_2}(T) \dots w_{x_m}(T)$$

the string is uniquely decodable (prefix free): just follow the game,



*or instantaneous, stronger than u.c., can be decoded on the fly reading only once the string*

if the  $x_i$  are iid with law  $p(x)$ , the length of the total sequence will be  $n \bar{L}(T)$

$H(p)$  is the average number of bits (within 1) per symbol necessary to compress a sequence of symbols generated by a source of prob  $p$

rk: use gzip to measure the entropy of English

- not iid, but short correlations
- this is for loss-less compression, in image / sound compression some mistakes are tolerable, compromise between accuracy of reconstruction and rate of compression: rate-distortion theory

### II.D. Mutual information

useful for the following and for other lectures (Florent)

$(X, Y)$  a pair of (a priori correlated) r.v. on  $(X \times X')$  not necessarily the same

$$P_{X,Y}(x,y) = \mathbb{P}[X=x \text{ and } Y=y] \quad \text{joint law}$$

$$\left. \begin{aligned} P_X(x) &= \sum_{y \in X'} P_{X,Y}(x,y) \\ P_Y(y) &= \sum_{x \in X} P_{X,Y}(x,y) \end{aligned} \right\} \text{marginal laws}$$

$$P_{X|Y}(x|y) = \frac{P_{X,Y}(x,y)}{P_Y(y)} \quad \text{conditional law, } \sum_x P_{X|Y}(x|y) = 1$$

if  $X, Y$  ind,  $P_{X,Y}(x,y) = P_X(x) P_Y(y)$ ,  $P_{X|Y}(x|y) = P_X(x)$

various entropies:  $H(X, Y) = - \sum_{x,y} P_{X,Y}(x,y) \log_2 P_{X,Y}(x,y)$

$$H(X) = - \sum_x P_X(x) \log_2 P_X(x)$$

$$H(Y) = - \sum_y P_Y(y) \log_2 P_Y(y)$$

$$\begin{aligned} H(X|Y) &= \sum_y P_Y(y) \left( - \sum_x P_{X|Y}(x|y) \log_2 P_{X|Y}(x|y) \right) \\ &= - \sum_{x,y} P_{X,Y}(x,y) \log_2 P_{X|Y}(x|y) \end{aligned}$$

mutual information between  $X$  and  $Y$ :

$$I(X; Y) = D(P_{X,Y} \parallel P_X P_Y) = \sum_{x,y} P_{X,Y}(x,y) \log_2 \frac{P_{X,Y}(x,y)}{P_X(x) P_Y(y)}$$

properties:  $I \geq 0$  (we have seen this  $\forall D$ )

$$I = 0 \iff P_{X,Y} = P_X P_Y \iff X \text{ and } Y \text{ independent}$$

$$\begin{aligned} I(X; Y) &= \sum_{x,y} P_{X,Y}(x,y) \log_2 \left( \frac{P_{X|Y}(x|y)}{P_X(x)} \right) = H(X) - H(X|Y) \\ &= H(Y) - H(Y|X) \\ &= H(X) + H(Y) - H(X, Y) \end{aligned}$$

hence  $H(X|Y) \leq H(X)$  conditioning reduces entropy

$I(X; Y)$  measures how much you know (in bits) about one of the two r.v. if the other is revealed to you

also,  $H(X, Y) \leq H(X) + H(Y)$

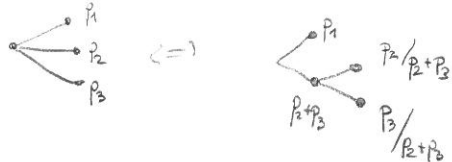
additional remarks, exercises:

- proof that canonical Gibbs Boltzmann is the one maximizing entropy under constraint on average energy
- proof that uniform law is the one with maximal entropy
- in Shannon's paper, proof that if one assumes

$\rightarrow H(p_1, \dots, p_M)$  continuous

$\rightarrow H(\frac{1}{M}, \dots, \frac{1}{M}) \uparrow$  with  $M$

$\rightarrow H$  "additive" under decompositions

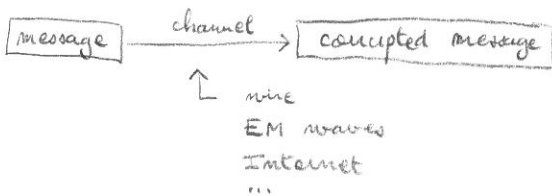


then only possibility is  $-\sum_i p_i \ln p_i$  within a multiplicative constant

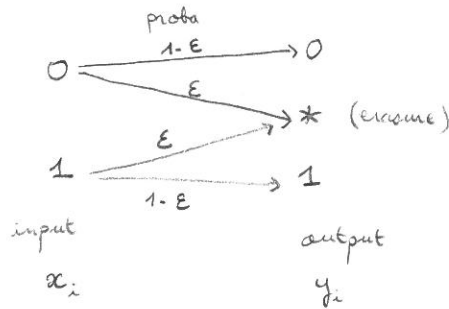


### III Communication over noisy channels

#### III.A Definitions



#### examples • Binary Erasure Channel (BEC)



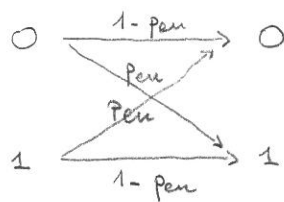
$E=0$  no noise

$E=1$  no signal at all

symmetric in 0/1

if one receives 0/1 one is sure that it was the correct value

#### • Binary Symmetric Channel (BSC)



$p_{err}=0$  no noise

$p_{err}=\frac{1}{2}$  no signal at all

restrict  $p \leq \frac{1}{2}$  by symmetry

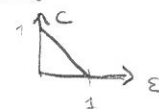
if  $p_{err} > 0$  one is never sure of the result

Capacity of a channel : X the input, Y the output,

$$C = \max_{P_X} I(X; Y) \quad : \text{nb of "effective bits of information" transmitted per use of the channel}$$

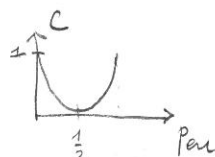
exercises: compute C for BEC and BSC, show that

•  $C = 1 - E$  for the BEC



•  $C = 1 - h(p_{err})$  for the BSC,

$$h(p) = -p \ln p - (1-p) \ln (1-p)$$



meaning will be clearer later on

proof of the capacity

BSC  $p_X$  parametrized by  $a, p_X(0)$

$$I(X; Y) = H(Y) - \underbrace{H(Y|X)}_{h(p)}$$

independently on  $a$

$\rightarrow Y=0$  with proba  $a(1-p) + (1-a)p$

$$\Rightarrow H(Y) = h(a(1-p) + (1-a)p) \text{ maximized with } a = \frac{1}{2}, H(Y) = h\left(\frac{1}{2}\right) = 1$$

obvious a posteriori by symmetry

$$C = 1 - h(p)$$

BEC  $H(Y|X) = h(\epsilon)$  independently on  $a$  either  $\{0, *\}$  or  $\{1, *\}$

$$H(Y): Y = \begin{cases} 0 & \text{with proba } a(1-\epsilon) \\ 1 & (1-a)(1-\epsilon) \\ * & \epsilon \end{cases}$$

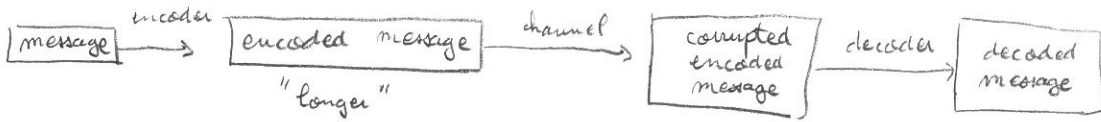
$$H(Y) = -\epsilon \ln \epsilon - a(1-\epsilon) \ln(a(1-\epsilon)) - (1-a)(1-\epsilon) \ln((1-a)(1-\epsilon))$$

$$= -\epsilon \ln \epsilon - (1-\epsilon) \ln(1-\epsilon) - (1-\epsilon) a \ln a - (1-\epsilon)(1-a) \ln(1-a)$$

$$\Rightarrow I(X; Y) = (1-\epsilon) h(a), \text{ max in } a = \frac{1}{2} \text{ (again obvious)} \Rightarrow C = 1 - \epsilon$$

### Encoding and decoding

noise of the channel destroys information  $\Rightarrow$  fight it by transmitting more information, i.e. add redundancy



sender and receiver agree on the encoding and decoding procedure beforehand

rate of a code :  $\frac{\text{nb of bits of message}}{\text{nb of bits of encoded message}} < 1$ , measure of redundancy, should be as big as possible

### III.B Naive coding

simplest way to be redundant: repeat oneself

encoding :  $0 \rightarrow 000$   
 $1 \rightarrow 111$

on the BSC, with  $p < \frac{1}{2}$

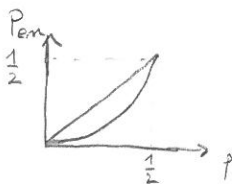
decoding :  $000 \rightarrow 0$   
 $001, 010, 100 \rightarrow 0$   
 $011, 101, 110 \rightarrow 1$   
 $111 \rightarrow 1$

majority rule (better take an odd number of repetitions to avoid it)

we want to transmit 0 or 1 with the same proba

without repetition,  $P_{\text{err}} = p$  (just take output as a guess)

with repetition,  $P_{\text{err}} = 3p^2(1-p) + p^3$  : two or three flips induce a mistake one can be corrected



better, but : \*  $R = \frac{1}{3}$

\*  $P_{\text{err}} > 0 \quad \forall p > 0$

to have  $P_{\text{err}} \rightarrow 0$  repeat 5, 7, ... times, but then  $R \rightarrow 0$   
can one do better ?

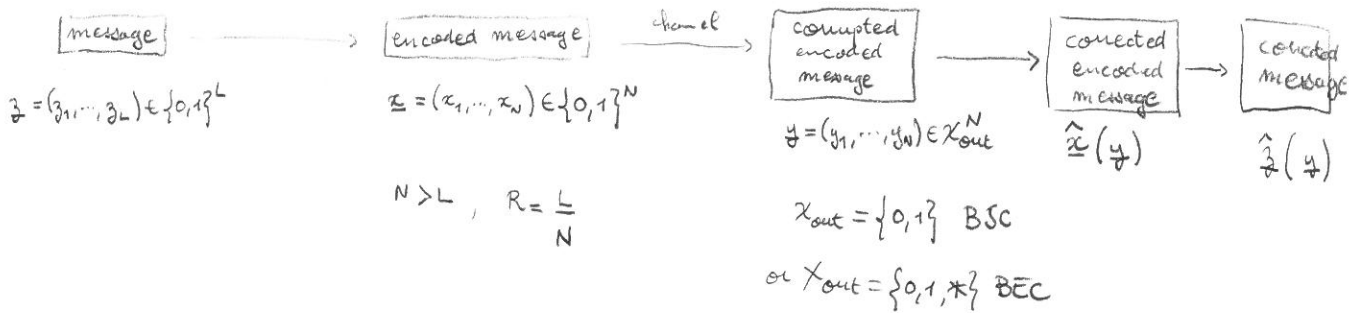
### III.C. Shannon channel coding theorem

quite surprisingly, answer of Shanna than is yes;

$\exists$  codes (of growing size) with  $P_{\text{err}} \rightarrow 0$  for any  $R < C$

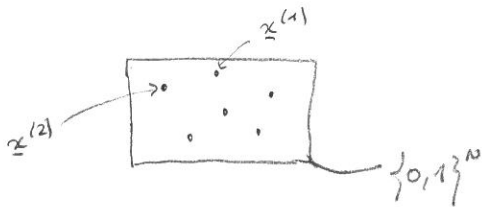
- \*  $C$  is the ultimate limit for the rate (Fano inequality) intuitive given the interpretation of  $I$
- \* yet it can be approached arbitrarily close
- \* statement true in the thermodynamic limit

#### more formal definitions



encoding:  $\underline{x} = f_{\text{encoding}}(\underline{z})$  injective

ie we choose  $2^L$  "codewords"  $\underline{x}^{(\alpha)} \in \{0,1\}^N, \alpha = 1, 2, \dots, 2^L$  to represent the possible messages



codebook:  $\mathcal{C} \subset \{0,1\}^N$   
 $\mathcal{C} = \{ \underline{x}^{(1)}, \underline{x}^{(2)}, \dots, \underline{x}^{(2^L)} \}$

intuitively, if  $\underline{x}^\alpha$  is transmitted,  $\underline{y}$  has moved because of the noise,

to avoid mis-recognition the  $\underline{x}^{(\alpha)}$  should be "far away" from each other  $\Rightarrow$  packing problem in Hamming space, cf Cohn's lectures  
(but geometry in high dimension counter-intuitive, 3 points are never "aligned" in the same plane as in 2d)

$\hat{\underline{x}}(\underline{y})$ : estimation of the sent  $\underline{x}$  given the received  $\underline{y}$

then  $\hat{\underline{z}}(\underline{y})$  is just  $f_{\text{encoding}}^{-1}(\hat{\underline{x}}(\underline{y}))$

assuming the  $2^L \underline{z}$  are equiprobable, then the  $2^L \underline{x}^\alpha$  are equiprobable, we can forget about  $\underline{z}$  and  $\hat{\underline{z}}$ : the problem is then:  
 $\underline{x} \in \mathcal{C}$  u.a.n,  $\rightarrow \underline{y}$  corrupted  $\rightarrow$  reconstruct  $\underline{x}$

decoding as an inference problem

(for one given codebook)

$\underline{x}$  random codeword,  $P_{\underline{x}}(\underline{x}) = \frac{1}{2^L} \mathbb{1}(\underline{x} \in \mathcal{B})$  indicator function

$\underline{y}$  random output  $P_{\underline{y}|\underline{x}}(\underline{y}|\underline{x}) = \prod_{i=1}^N Q(y_i|x_i)$

with  $Q$  describing the channel, for BEC

- $Q(0|0) = 1 - \epsilon$
- $Q(*|0) = \epsilon$
- $Q(1|1) = 1 - \epsilon$
- $Q(*|1) = \epsilon$
- $Q(1|0) = Q(0|1) = 0$

given  $\underline{y}$  we have to guess  $\underline{x}$

Bayes theorem:  $P_{\underline{x}|\underline{y}}(\underline{x}|\underline{y}) = P_{\underline{y}|\underline{x}}(\underline{y}|\underline{x}) \frac{P_{\underline{x}}(\underline{x})}{P_{\underline{y}}(\underline{y})} = \frac{1}{Z(\underline{y})} \mathbb{1}(\underline{x} \in \mathcal{B}) \prod_{i=1}^N Q(y_i|x_i)$

as there is noise we cannot a priori do better than assigning probabilities to the various possible input codewords

$\hat{\underline{x}}(\underline{y})$ ? depends on the measure of error one wants to minimize

$\hat{\underline{x}}(\underline{y}) = \underset{\underline{x}}{\operatorname{argmax}} P_{\underline{x}|\underline{y}}(\underline{x}|\underline{y})$  minimizes  $P[\hat{\underline{x}} \neq \underline{x}]$

block-Maximal A Posteriori (MAP) decoding  $\Leftrightarrow$  Maximum Likelihood (ML)

$\hat{x}_i(\underline{y}) = \underset{x_i}{\operatorname{argmax}} P_{x_i|\underline{y}}(x_i|\underline{y})$  minimizes the number of badly decoded bits

$\uparrow$  marginal law of  $P_{\underline{x}|\underline{y}}$

more on that in Krzakala's lectures

for the BEC:  $\underline{y} = (0, 1, *, *, 1, 0, *, 0)$

$\uparrow$  these are fixed

$\uparrow$  these are completely free

$P_{\underline{x}|\underline{y}}(\underline{x}|\underline{y})$ : uniform over those codewords which agree with  $\underline{y}$  on the uncoded place

call  $B(\underline{y}) = \{ \underline{x} \in \{0,1\}^N : \text{if } y_i \in \{0,1\} \text{ then } x_i = y_i \}$

$= \frac{1}{|B(\underline{y})|} \mathbb{1}(\underline{x} \in \mathcal{B} \cap B(\underline{y}))$

# Shannon Random Code Ensemble

reminiscent of REM (Biroli's lectures), also in the treatment

$\mathcal{C} = \{ \underline{x}^{(\alpha)} \}_{\alpha=1, \dots, 2^L}$  obtained by taking  $x_i^{(\alpha)} = \begin{cases} 0 & \text{proba } 1/2 \\ 1 & \text{proba } 1/2 \end{cases}$  independently  $\forall i, \alpha$

very simple, seems crazy, does not ensure that the  $\underline{x}^{(\alpha)}$  are all different, but in the limit  $N, L \rightarrow \infty$  very few collisions

on the BEC: suppose w.l.o.g. that  $\underline{x}^{(1)}$  has been transmitted

$N_*(\underline{y})$ : nb of erased bits in  $\underline{y}$ , r.v.  $\text{Bin}(N, \epsilon) \simeq N\epsilon$  in the thermodynamic limit

how many codewords  $\neq \underline{x}^{(1)}$  in  $\mathcal{B}(\underline{y})$ ?

$$\text{i.o. } \mathcal{N} = \text{Bin} \left( 2^L - 1, \left( \frac{1}{2} \right)^{N - N_*(\underline{y})} \right)$$

$\leftarrow$  proba that  $\underline{x}^{(\alpha+1)}$  agrees with  $\underline{x}^{(1)}$  on the non erased bits

$$\mathbb{E}[\mathcal{N}] \simeq 2^{L-N+\epsilon N} = 2^{N(R-1+\epsilon)} \rightarrow 0 \text{ if } R < 1 - \epsilon = C \text{ as } N \rightarrow \infty$$

$\Rightarrow$  as in the REM, a first moment method,  $\mathcal{N} = 0$  w.h.p

$\Rightarrow$  one can recover  $\underline{x}^{(1)}$  as the only codeword compatible with the received bits in  $\underline{y}$

$\Rightarrow$  proves Shannon's channel coding theorem for the BEC,

one can decode with vanishing error probability with rates up to capacity

proof for the BSC as homework, text and solution (in french) on the webpage,

⚠ notations inverted

very simple random code achieves capacity

BUT: needs exponentially large (in  $N$ ) memory to store the codewords

needs exponential time to decode

$\Rightarrow$  needs to add structure to the codebook, can this be done and still achieves capacity? to be seen in the following

rk: here we have average <sup>with respect / codewords and to the code</sup> error probability  $\rightarrow 0$ , can be boosted to maximal (over the codewords) error probability  $\rightarrow 0$ , hence  $\exists$  one code with max error proba  $\rightarrow 0$  by expurgating the worst half codewords, does not change the rate

- converse with Fano,  $P_{\text{er}} > 0$  if  $R > C$ , actually  $P_{\text{er}} \rightarrow 1$  (and limit in  $N$ )

# IV Low Density Parity Check Codes (LDPC)

## IV.A. Linear codes

need to add some structure to the codebook  $\mathcal{C} = \{x^{(c)}\} \subset \{0,1\}^N$  to make encoding and decoding easier

$\{0,1\}^N$  is a linear space over  $\{0,1\} = \mathbb{Z}_2$ , with

$0+0=0$	$0 \cdot 0 = 0$
$0+1=1+0=1$	$0 \cdot 1 = 0$
$1+1=0$	$1 \cdot 0 = 0$
addition mod 2	$1 \cdot 1 = 1$

$\mathcal{C}$  is a linear code if it is a linear subspace of  $\{0,1\}^N$   
 ie  $\left\{ \begin{array}{l} 0 \in \mathcal{C} \text{ and} \\ x, y \in \mathcal{C} \end{array} \right. \Rightarrow x+y = (x_1+y_1, \dots, x_N+y_N) \in \mathcal{C}$   
↑ mod 2    ↑

would be enough,  $x+x=0 \Rightarrow$  imply first

linear algebra on  $\{0,1\}$  instead of  $\mathbb{R}$ , works pretty the same

$\dim \mathcal{C}$ : nb of linearly independent codewords  $\Rightarrow |\mathcal{C}| = 2^{\dim \mathcal{C}}$

$\mathcal{C}$  can be specified as  $\text{Ker } H = \{x \in \{0,1\}^N, Hx = 0\}$ ,  $H$  a  $M \times N$  matrix on  $\mathbb{Z}_2 \Rightarrow 0,1$ -wise

or as  $\text{Im } G = \{uG, u \in \{0,1\}^{N-M}\}$ ,  $G$  a  $(N-M) \times N$  matrix

$H$  parity check matrix, one line of  $Hx=0$  of the form  $x_{i_1} + x_{i_2} + x_{i_3} = 0$   
 parity of the number of 1's must be even

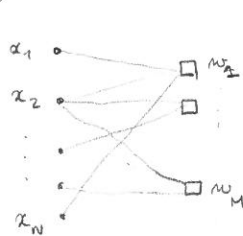
if these equations are linearly independent,  $|\mathcal{C}| = 2^{N-M}$

$\dim \text{Ker } H = N - \text{rk } H$  always, independence  $\Leftrightarrow$  full rank

rate of the code  $R = \frac{N-M}{N} = 1 - \frac{M}{N}$  (assuming equations independent, in particular  $M < N$ )

$G$  generator matrix, allows to encode the codewords very easily

Tanner graph representation of parity checks:



$$m_a(x) = 1 \left( \sum_{i=1}^N H_{ai} x_i = 0 \right)$$

edge between  $a$  and  $i \Leftrightarrow H_{ai} = 1$

$\partial a = \{i : H_{ai} = 1\}$ ,  $\partial i = \{a : H_{ai} = 1\}$  neighborhoods on the graphs

all codewords are equivalent in a linear code, and  $(0, \dots, 0)$  is always a codeword

$\mathbb{F}_2$ : addition mod 2 of  $\{0,1\} \Leftrightarrow$  exclusive OR of  $\{\text{True}, \text{False}\} \rightarrow$  known as a CSP as XORSAT

$\sigma_i = (-1)^{x_i} = \begin{cases} 1 & x_i = 0 \\ -1 & x_i = 1 \end{cases}$   $\sum_i x_i = 0 \text{ mod } 2 \Leftrightarrow \prod_i \sigma_i = 1$  : also p-spin interactions  
↑ mentioned per this  
 $\Rightarrow$  lectures of Biale, Ricci-Tersenghi, Messia...

### IV B. Definition of the simplest ensemble

given the success of Shannon's random code, try a random  $H$  (Ballager 62)

but "low density" : a few 1's by row and column of  $H \rightarrow$  few operations, faster computationally

take  $H$  such that :  $|D_i| = l \quad \forall i$  , i.e.  $l$  1's in each column  
 $|D_a| = k \quad \forall a$  ,  $k$  rows

with  $l, k$  finite in the thermodynamic limit

$\Rightarrow Nl = Mk$  , total number of 1's in the matrix  $R = 1 - \frac{l}{k}$

how to do this in practice:  $l < k$



random matching (permutation) of these  $Nl = Mk$  half edges

$l=3$   
 $k=4$   
Here



there can be pb if



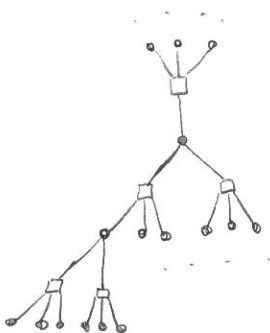
can eliminate even nb of parallel edges

$$\alpha_i + \alpha_i = 0 \pmod{2} \quad \forall \alpha_i$$

does not really matter, will be rare in the thermodynamic limit

if  $M < N$  the equations are linearly independent (or sub-extensive nb are dependent)

crucial property : locally tree-like in the thermodynamic limit



no short loops (with high probability)

intuitive explanation: exploring from one vertex, choose finite nb of neighbors out of an extensive one, proba to choose twice in the same finite set is  $O(1/N)$

$\Delta$  not the tree of the guessing game



# IV. C Analysis on the BEC

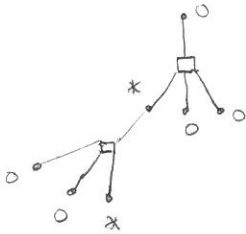
assume codeword  $(0, \dots, 0)$  has been transmitted over BEC( $\epsilon$ )

→ because all codewords are equivalent, but we must pretend not to know that it is  $(0, \dots, 0)$

received  $y = (0, *, 0, 0, \dots, *)$

sender and receiver have agreed on an LDPC

what can the receiver do?



the first \* can be set to 0 with the first parity check

then the second \* as well

- in all checks where there is a single \* → set it to 0
- continue iteratively

when it stops, either no \* remains → perfect decoding

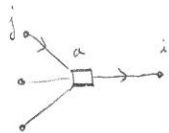
• or at least two \* in the checks that remain

stopping set: subset of the var such that each check contains either 0 or ≥ 2 var of the stopping set

polynomial time algorithm

up to which value of  $\epsilon$  will the first situation occur?

reformulation of the algorithm as message passing (cf Moore's lectures)

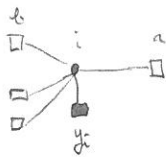


$$u_{a \rightarrow i} = \begin{cases} 0 & \text{if } h_{j \rightarrow a} = 0 \quad \forall j \in \partial a \setminus i \\ * & \text{otherwise} \end{cases}$$

$$u_{a \rightarrow i}, h_{i \rightarrow a} \in \{0, *\}$$

"I'm sure you are 0"

"I can't say"



$$h_{i \rightarrow a} = \begin{cases} * & \text{if } y_i = * \text{ and } u_{b \rightarrow i} = * \quad \forall b \in \partial i \setminus a \\ 0 & \text{otherwise} \end{cases}$$

"I'm sure I'm a 0, either because I've received my bit, or deduced it from one neighboring check"

$$h_i = \begin{cases} * & \text{if } y_i = * \text{ and } u_{a \rightarrow i} = * \quad \forall a \in \partial i \\ 0 & \text{otherwise} \end{cases}$$

final estimate

in discrete time  $u_{a \rightarrow i}^{(t=0)} = *$  for all  $a \rightarrow i$

$$h_{i \rightarrow a}^{(t=0)} = f(\{u^{(t=0)}, y\})$$

$$u_{a \rightarrow i}^{(t=1)} = f(\{h^{(t=0)}\})$$

$h_i^{(t)}$ : estimate for the  $u^{(t)}$

monotonicity, messages switch from \* to 0, at the end same final state for the  $h_i$  than the original decoder, for any update schedule

probabilistic analysis

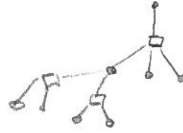
$$y^{(t)} = \text{Proba} (u_{a \rightarrow i}^{(t)} = 0)$$

with respect to

- choice of the H
- choice of the pattern of erasures (channel)
- choice of an edge uniformly at random

$$z^{(t)} = \text{Proba} (h_{i \rightarrow a}^{(t)} = 0)$$

idem



$$\left\{ \begin{aligned} z^{(t=0)} &= 0 \\ z^{(t+1)} &= (z^{(t)})^{k-1} \end{aligned} \right.$$

all  $i$ -puls must be 0, id because of the tree structure

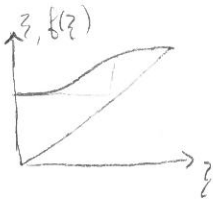
$$z^{(t)} = 1 - \epsilon (1 - z^{(t-1)})^{k-1}$$

all must be \* for  $h$  to be \*

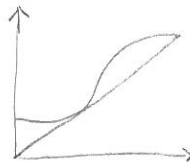
$$z^{(t)} = f(z^{(t-1)})$$

$$f(z) = 1 - \epsilon (1 - z^{k-1})^{k-1}$$

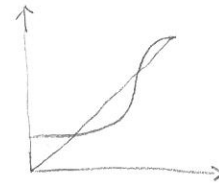
properties of  $f$ :  $f(0) = 1 - \epsilon$  if  $k \geq 3, \epsilon \geq 3$ , zero derivatives in 0 and 1  
 $f(1) = 1$   
 $f' \uparrow$



$\epsilon$  small  
 $\epsilon < \epsilon_{BP}$



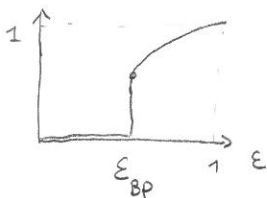
$\epsilon = \epsilon_{BP}$



$\epsilon$  big  
 $\epsilon > \epsilon_{BP}$

plot of the iteration, plateau in  $z^{(t)}$ , exponent 1/2, similarity with  $C(t)$

$$\mathbb{P} [h_{i \rightarrow a}^{(t=\infty)} = *] = \epsilon (1 - z^{(\infty) k-1})^{\epsilon}$$



for  $E < E_{BP}$ , perfect decoding, in linear time, for a positive rate  $\Rightarrow$  non trivial results

numerical values

$l$	$k$	$E_{BP}$	$E_{MAP}$	$R$	$E_{SR}$
3	4	0,647	0,746	$1/4 = 0,25$	0,75
3	5	0,518	0,591	$2/5 = 0,4$	0,6
3	6	0,429	0,488	$1/2 = 0,5$	0,5
4	6	0,506	0,665	$1/3 = 0,333$	0,666

not yet defined

$E_{SR} = 1 - R$  the maximal level of noise a code of this rate could correct

$E_{BP} < E_{SR}$  is this because of the code ?

or because of the algorithm to decode ?

in other words, for  $E > E_{BP}$ , once BP has inferred all possible implications,

one has a linear system of equations on  $N'$  variables,  $M'$  equations,

has this system a single solution ? then perfect decoding would still be possible by exhaustive algorithm (or Gauss elimination here)

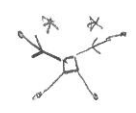
or more than one solution, then perfect decoding is not possible, whatever the algorithm

$$N' = N \epsilon (1 - \zeta^{k-1})^\ell$$

with  $\zeta = \zeta^\infty$ , the number of undecoded variables

$$M' = M \left( 1 - \zeta^k - k \zeta^{k-1} (1 - \zeta) \right)$$

at least two in the incoming messages



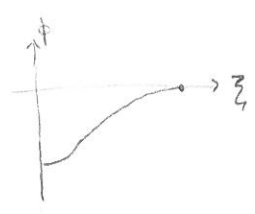
one needs to work with messages and not  $h_i$  to have independence

assuming the equations are independent,

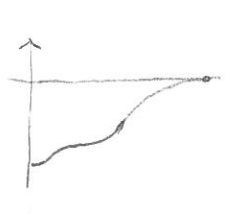
$$\text{nb of solutions} = 2^{N'-M'} = 2^{N\phi}$$

$$\text{with } \phi(\epsilon, \zeta) = \epsilon (1 - \zeta^{k-1})^\ell - \frac{\ell}{k} \left( 1 - \zeta^k - k \zeta^{k-1} (1 - \zeta) \right)$$

$$\begin{aligned} \frac{d}{d\zeta} \phi &= -\ell \epsilon (k-1) \zeta^{k-2} (1 - \zeta^{k-1})^{\ell-1} + \ell \left( \zeta^{k-1} + (k-1) \zeta^{k-2} - k \zeta^{k-1} \right) \\ &= \ell (k-1) \zeta^{k-2} \left[ -\epsilon (1 - \zeta^{k-1})^{\ell-1} + 1 - \zeta \right] \\ &= 0 \text{ when } \zeta \text{ is a fixed point of the iterations} \end{aligned}$$

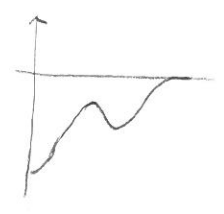


$\epsilon < \epsilon_{BP}$

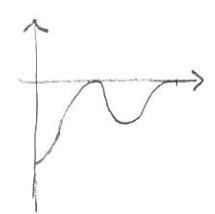


$\epsilon = \epsilon_{BP}$

inflection point



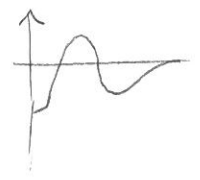
$\epsilon \in [\epsilon_{BP}, \epsilon_{MAP}]$



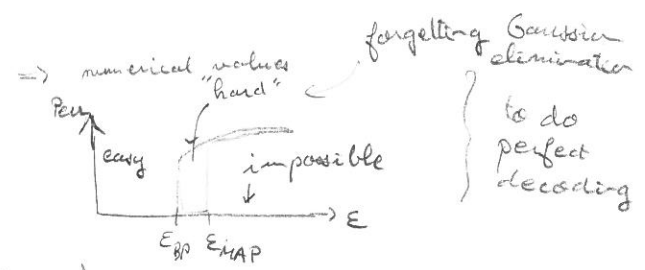
$\epsilon = \epsilon_{MAP}$

like Franz-Pearce potential (with - sign,  $1 - \zeta \leftrightarrow q$ )

plots shown by Florent  
but not exactly the same phases,  $m=0$  vs  $m=1$  in terms of magnetization



$\epsilon > \epsilon_{MAP}$



computational gap, as in many inference pb (Florent)  
 $\approx$  like  $T_d / T_c$ , but not completely

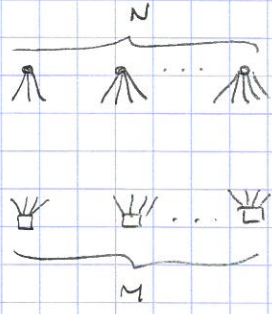
IV.D. More generic LDPC ensembles

instead of fixed degrees  $(l, k)$  distribution of degrees  
variables checks

$\lambda_e$ : fraction of variables that have degree  $e$

$\rho_k$ : checks  $k$

H uniform under this constraint



← permutation

$\langle e \rangle = \sum_e e \lambda_e$ ,  $\langle k \rangle = \sum_k k \rho_k$   $N \langle e \rangle = M \langle k \rangle$  to match

still locally tree like

if one draws at random an edge, ends up on  $\circ$  with degree  $l+1$

with proba  $\tilde{\lambda}_e = \frac{(l+1) \lambda_{e+1}}{\langle e \rangle}$  on  $\square$  with degree  $k+1$

with proba  $\tilde{\rho}_k = \frac{(k+1) \rho_{k+1}}{\langle k \rangle}$

define  $\lambda(x) = \sum_e \lambda_e x^e$ ,  $\rho(x) = \sum_k \rho_k x^k$   
 $\tilde{\lambda}(x) = \sum_e \tilde{\lambda}_e x^e$ ,  $\tilde{\rho}(x) = \sum_k \tilde{\rho}_k x^k$   $\tilde{\lambda}(x) = \frac{\lambda'(x)}{\lambda'(1)}$

then, exercise:  $\begin{cases} \eta^{(t+1)} = \tilde{\rho}(\zeta^{(t)}) \\ \zeta^{(t)} = 1 - \epsilon \tilde{\lambda}(1 - \eta^{(t)}) \end{cases}$

same reasoning as before, just add the proba of finding a vertex of a certain degree

$\Phi(\epsilon, \zeta) = \epsilon \lambda(1 - \tilde{\rho}(\zeta)) - \frac{\langle p \rangle}{\langle k \rangle} \left( 1 - \rho(\zeta) - \langle k \rangle (1 - \zeta) \tilde{\rho}(\zeta) \right)$

generalizes the case with fixed  $k, l$

one can find choices of  $\{\lambda_e, \rho_k\}$  such that  $E_{BP}$  is arbitrarily close to  $E_{Sh}$

seems great but  
linear decoding up to capacity

not obvious that they are the best in practice, depends on the finite  $N$  behavior, and prefactor in the linear complexity diverges when  $E_{BP} \rightarrow E_{Sh}$  (max. degree  $\uparrow$ )

### V. Conclusions and perspectives

what we have done is a very specific (and simple) example of the reasoning of the cavity method: Let's abstract the main points, besides the IT context:

- variables  $\underline{x} = (x_1, \dots, x_N) \in \mathcal{X}^N$  (here the encoded message to be recovered)
- quenched randomness  $\underline{J}$  (here the graph/parity check matrix + realization of the noise)
- defining a factor graph with  $M$  interactions

• probability law 
$$\mu(\underline{x}; \underline{J}) = \frac{1}{Z(\underline{J})} \prod_{i=1}^N w_i(x_i; \underline{J}) \prod_{a=1}^M w_a(\underline{x}_{\partial a}; \underline{J})$$

$\uparrow$  here: information received from the channel       $\uparrow$  constraints on the codewords

goals of the cavity method:

- compute  $\phi = \lim_{N \rightarrow \infty} \frac{1}{N} \mathbb{E} [\ln Z(\underline{J})]$  quenched free entropy density (here  $\simeq$  conditional entropy)
- marginal laws  $\mu(x_i; \underline{J})$  "magnetizations"

strategy: tree like factor graph

- trees exactly solved by message passing (here the messages were 0/1 in general no trivial probability distribution on  $\mathcal{X}$ )
- long loops  $\rightarrow$  "irrelevant"  $\rightarrow$  RS cavity method  $\rightarrow$  induce long range correlations  $\rightarrow$  RSB

⚠ LDPC = planted models,  $\simeq$  the Mishchenko line, don't take as completely generic the phenomena / transitions found here

$\Rightarrow$  cf Federico's lectures

Scemas 89  
 Richardson - Urbanke } 2001  
 Wiboy - Mitzemacher - Shokrollahi - Spielma  
 Kabashima - Saad  $\simeq$  2000  
 Montanari  $\simeq$  2000