

Chapter 2 : Compression / Prediction / Generation equivalence

I Concepts

- Lossless compression

data = sequence of symbols, from a known predefined alphabet (discrete setting)

↳ set of symbols

ex: $\{ 'a', 'b', 'c', \dots, 'z' \}$

ex: $\{ '1', '2', '3', \dots \}$

ex: $\{ 'a', 'p', 'q', \dots \}$

- Lossless compression:

↳ no information lost

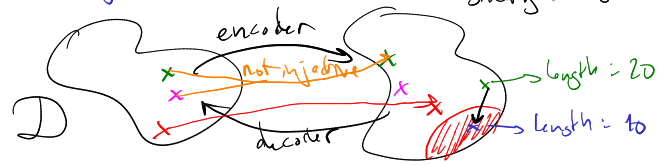
encoder: data \mapsto binary string

⚠ injective
otherwise: not decodable
↳ "uniquely decodable code"

decoder: binary string \mapsto data

if encoder not surjective: not optimal compression
Binary strings

lossless \Leftrightarrow decoder \circ encoder = Id



- Prediction

given a sequence of observed variables: x_1, x_2, \dots, x_n , predict x_{n+1}

i.e. give the probability distribution of each possible symbol (discrete setting \rightarrow finite alphabet)
ex: $\{ 'a', 'b', 'c' \}$

predictor: function: data $(x_1, \dots, x_n) \mapsto p(x_{n+1} | x_1, \dots, x_n)$ probability of next variable

key: vector: $\begin{pmatrix} p('a' | x_1, \dots, x_n) \\ p('b' | x_1, \dots, x_n) \\ p('c' | x_1, \dots, x_n) \end{pmatrix}$
 $\Sigma = 1$

probability distribution over the next variable x_{n+1}

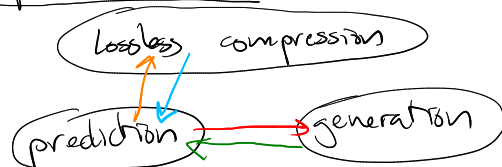
- Generative model

model able to generate new data

data D is generated according to its model probability $p(D)$ is often as

2 points \rightarrow distribution over data
 \rightarrow sample from it

II From one concept to the other



- Prediction \rightarrow generation [easy]

start from $x_1 \rightarrow$ given
 \rightarrow or taken from learned distrib

choose x_{n+1} according to the probability law $p(x_{n+1} | x_1, \dots, x_n)$ used to predict next symbol

iterate \rightarrow generate data this way

$$p(D) = \prod_i p(x_i | x_1, \dots, x_{i-1})$$

↳ x_1, \dots, x_n

$$= p(x_2) \times p(x_2 | x_1) \times p(x_3 | x_1, x_2) \times p(x_4 | x_1, \dots, x_3) \dots$$

$$= p(D) \quad \text{Bayes theorem applied } n-1 \text{ times}$$

Generation \rightarrow Prediction [By sampling... might be slow]

- given $x_1 \dots x_n$, how to estimate the probability distribution of x_{n+2} ?
- generate new data many many times, according to the generator law
- check how frequent each $(x_1 \dots x_n, x_{n+2})$ is generated $\rightarrow p(x_1, \dots, x_{n+2})$
- compute $p(x_{n+2} | x_1 \dots x_n) = p(x_1, \dots, x_{n+2}) / p(x_1, \dots, x_n)$ Bayes

Compression \rightarrow Prediction

- idea: strings whose encoding is shorter are more probable
 \hookrightarrow coherent with $L = -\log p \rightarrow p = 2^{-L}$

\hookrightarrow complete a string:
 $L(\text{string} + 'c') \rightarrow p('c')$ or $2^{-L(\text{string} + 'c')}$

- given $x_1 \dots x_n$, build all possible completed sequences $x_1 \dots x_n, x_{n+1}$
 sequence $s_k: x_1 \dots x_n, k^{\text{th}}$ symbol

- compress each of them independently, and check their encoding length $L_k = \text{size}(\text{encode}(s_k))$

- define $p(k) = 2^{-L_k}$

- \uparrow p needs to be a probab. distrib: $\sum_k p(k) = 1$

\hookrightarrow normalize $(p(k))$ by a global factor if necessary
 $k \in \text{alphabet}$

\hookrightarrow we'll see later that the norm is already 1 if the encoder is optimal

- define $p(x_{n+2} = k | x_1 \dots x_n) = p(k)$

ex: $A = \{a, b, c, \dots\}$

$x_1 \dots x_n =$ "This course is about"
 $p(x_{n+2} | x_1 \dots x_n)$? \hookrightarrow This course is about a \rightarrow 1001...
 b \rightarrow 001...
 c \rightarrow ...
 ...
 This course is about ...

Prediction \rightarrow Compression [the main point of this lesson]

- Lesson from entropy: if we have good prior information on the distribution $p(x_{n+2} | x_1 \dots x_n)$, then we don't need too much information to be communicated to know which symbol is chosen for x_{n+2}

- encode just that information \rightarrow so that it's decodable later

- entropy = average length of that code $E[-\log p] = H(p)$
 $x_{n+2} \sim p$

\hookrightarrow not possible to be shorter, without losing information

Can we prove that bound?
 \hookrightarrow is it possible to reach that bound? how?

III Encoders

- issue: when to stop reading data

* self-delimiting code

- \rightarrow : to code an integer n (not bounded)
 $n \in \mathbb{N}$

\hookrightarrow need $b = \lceil \log n \rceil$

\hookrightarrow need to tell the length of the sequence to decode

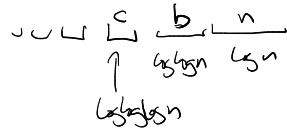
\hookrightarrow encode b : $\log b$ bits = $\log \log n$

1011100011|100
 one sequence to decode another one

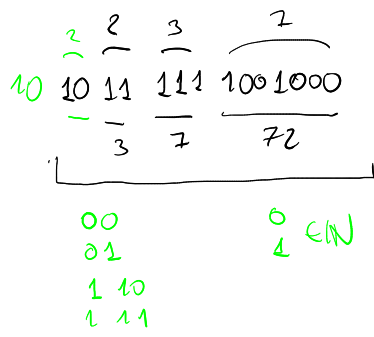
rewrite in binary

\uparrow
 \mathbb{N} 2 \rightarrow 10
 3 \rightarrow 11

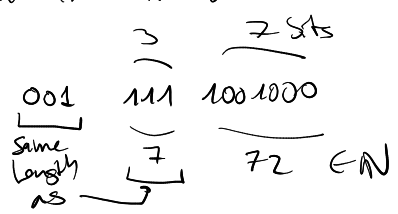
72 \rightarrow 64 + 8 = 1001000 $\lceil \log_2 n \rceil$



recursively until $\log \log \dots \log n \approx 1$
 total encoding length:
 $\log n + \log \log n + \log \log \log n + \dots$
 $+ \log \log \dots \log n$
 $= 2$



or stop at first iteration and encode differently



$$\lceil \log n \rceil + \lceil \log \log n \rceil \times 2$$

$$\log \log n \ll \log n$$

0000...01
 padding with 0
 flag $\log n - 1$ '0' followed by '1'

* codes with "end-of-file" character

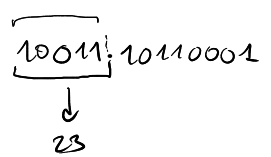
- postfix: read the data, until you meet a special code meaning "end"
- to do this: make binary blocks of fixed size and use one of them for "end"



- obvious cost: $-\log p(\text{"end"}) = 4$ + 4 bits
 - hidden cost: can't use symbols forming "end" at other places
 - cost: number of blocks $\times -\log(1 - p(\text{"end"}))$
 - 4-bit block: $2^4 = 16$ possible patterns
 $b-1$ for end
 15 remaining
- Factor ≈ 0.093 here
 \hookrightarrow not negligible if encoding long sequences

* Prefix codes

when the code of a string cannot be the beginning of the code of another string

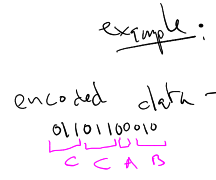
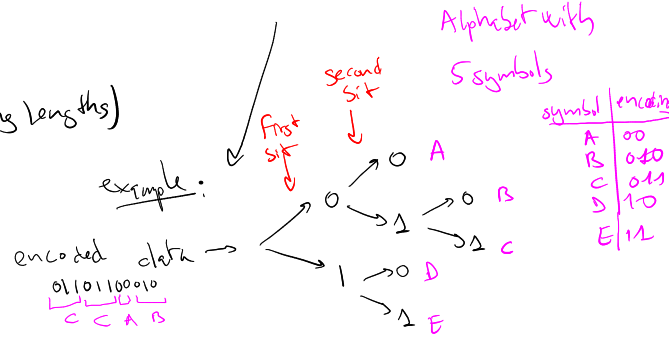


avoids confusion when decoding

David MacKay's book page 104 for examples

\hookrightarrow such codes can be written as binary trees

\hookrightarrow can show later, that without losing generality, one can assume that a code is prefix. Otherwise, rebuild another code with same coding lengths



IV Bounds, optimal codes and Kraft inequality

- entropy \Rightarrow make use of symbol probabilities (to produce shorter average encoding length)
 - \hookrightarrow more frequent symbols: should have shorter encoding lengths / code
- hint: $2^{-\text{length}(\text{code of } s)}$ should be $\propto p(s)$
 - \hookrightarrow expect $\text{Length}(\text{code of } s) = -\log p(s)$ \leftarrow cf. entropy

codes of each symbol in the alphabet

A: 001 \rightarrow 3
B: 101 \rightarrow 3

last ex.
 $L = \{2, 3, 3, 9, 2\}$

Kraft inequality

Given an encoder, let's denote by (l_i) the list of the lengths of all possible codewords (with multiplicity)

then:

the code is uniquely decodable $\Rightarrow \sum_i 2^{-l_i} \leq 1$

and reciprocally:

if (l_i) satisfies $\sum_i 2^{-l_i} \leq 1$, then there exists an encoder (prefix-code) whose codewords have this length distribution

Moreover: if equality ($\sum_i 2^{-l_i} = 1$), the code is said to be "complete"

Math proof: note $S = \sum_i 2^{-l_i}$

- choose any $n \in \mathbb{N}$

sketch

compute $S^n \rightarrow (\sum_i 2^{-l_i})^n \rightarrow$ count terms \rightarrow show it grows at most linearly with n
 $S^n \leq Kn$

$\begin{cases} l_{\min} = \min_i l_i \\ l_{\max} = \max_i l_i \end{cases} \quad \forall i, l_{\min} \leq l_i \leq l_{\max}$

$$S^n = \left(\sum_i 2^{-l_i} \right)^n = \left(\sum_i 2^{-l_i} \right) \times \left(\sum_j 2^{-l_j} \right) \times \left(\sum_k 2^{-l_k} \right) \dots$$

$$= \sum_{ijk \dots} \underbrace{2^{-l_i} \times 2^{-l_j} \times 2^{-l_k} \times \dots}_{= 2^{-(l_i + l_j + l_k + \dots)}}$$

$$= \sum_{l=n l_{\min}}^{n l_{\max}} A_l 2^{-l}$$

numbers of terms such that $l_1 + l_2 + \dots = l$

$\begin{matrix} 001 & 10 & 1 & 110 \\ i & j & k & \end{matrix}$ length $\Rightarrow l$

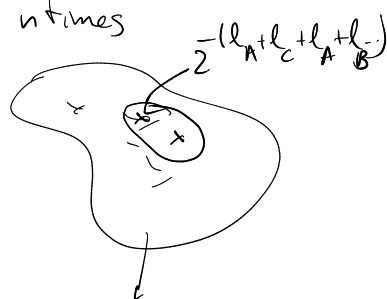
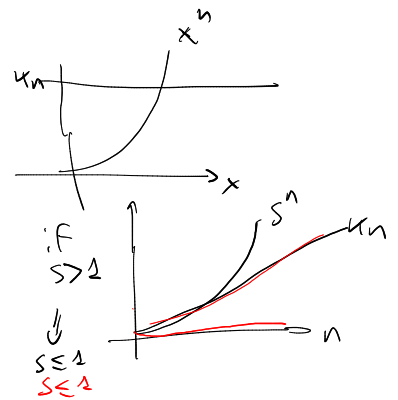
$$A_l \leq 2^l$$

how many terms

Maximum number of binary chains of length l

$$S^n \leq \sum_{l=n l_{\min}}^{n l_{\max}} 2^l \times 2^{-l} \leq n(l_{\max} - l_{\min}) \leq n l_{\max}$$

ex: $\overbrace{01100 \dots}^l 2 \rightarrow 2 \times 2 \dots = 2^l$



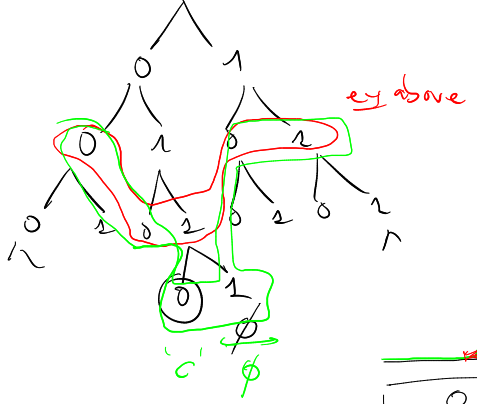
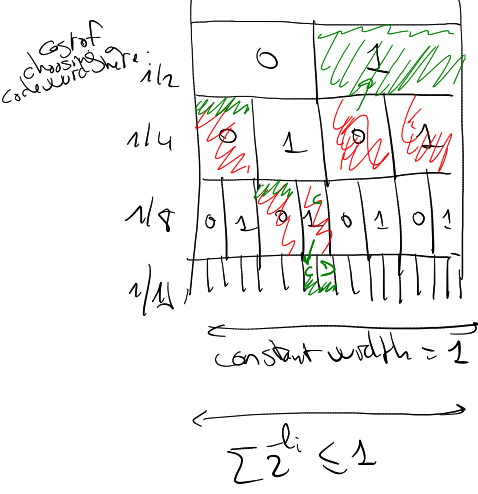
-0	-1	-2	-3	-4 ?
2	2	2	2	2
↓	↓	↓	↓	↓

Proof by drawing:

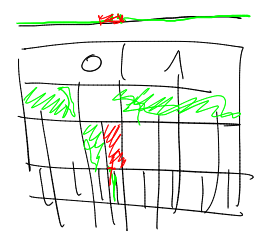
- list all the codewords

	A	B	C	D	E
	00	010	011	10	11
	00	01	10	11	?
ex:	0	1	2	2	?

\rightarrow not possible to have too many short codewords



For every code word taken, consider 2^{-l_i}



Lower bound on compression length

Length (encoding X) $\geq H(X)$

with equality only if the code length is coherent with the probability
 l_i (2^{-l_i}) $P(i)$ in the data
 \downarrow \downarrow
 $\log p(i)$

Proof: Gibbs inequality + Kraft

$$E[l(x)] = \sum_i p(i) l_i$$

$x \sim p$ distribution over symbols of the alphabet

$$= \sum_i p(i) (-\log S - \log q_i)$$

$q_i = \frac{2^{-l_i}}{\sum_j 2^{-l_j}}$ $i \rightarrow l_i \rightarrow q_i \propto 2^{-l_i}$

$$= -\log S - \sum_i p(i) \log q_i$$

$S \leq 1$ $\log S \leq 0$ $-\log S \geq 0$

$\sum q_i = 1$ q_i : probability

$q_i = \frac{2^{-l_i}}{S}$ $S q_i = 2^{-l_i}$

$l_i = -\log S - \log q_i$

\hookrightarrow Gibbs: cross entropy \geq entropy $-\sum p(i) \log q_i \geq H(p)$

$$\geq H(p)$$

Upper bound

It's possible to build a code such that:

Source coding theorem (symbol codes): There exists a variable-length encoding C of an ensemble X such that the average length of an encoded symbol, $L(C, X)$, satisfies:

$$L(C, X) \in [H(X), H(X) + 1]$$

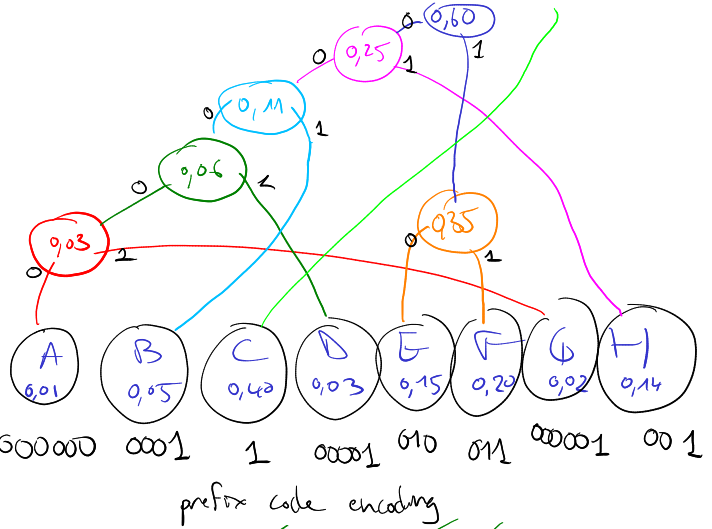
Proof: Kraft with $l_i = \lceil -\log p_i \rceil$

Constructive proof: \rightarrow Huffman coding (best possible with integer nb of bits/symbol) ex: 'A': 1010
 \rightarrow arithmetic " (non-integer nb of bits)
 \hookrightarrow reach $H(x)$ rate

V Huffman & arithmetic coding

Huffman coding

- 1) write down one node per symbol, with its probability
- 2) iteratively, join the 2 nodes with lowest probability, and form a new node, with prob = sum of the 2
- 3) this forms a binary tree
- 4) code = position in the tree



ex: BACH → 0001000001001

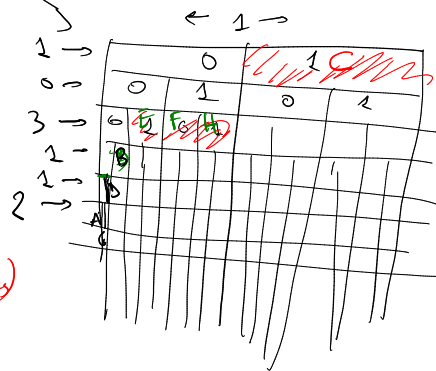
Optimality proof:

- frequencies match encoding lengths
- 1) for a given set of code word lengths,
 - 2) optimal set of code word lengths

Suppose another prefix code is better with a \neq set of code lengths

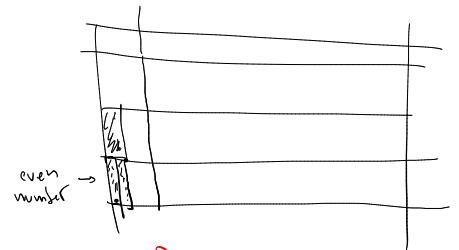
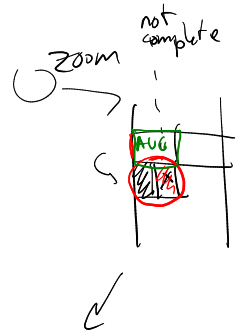
S_i : code lengths for least probable symbols are \neq

more true
 ⇒ length (2 least frequent symbols) = the same



$\{l_i\}$

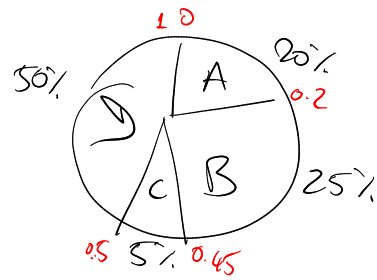
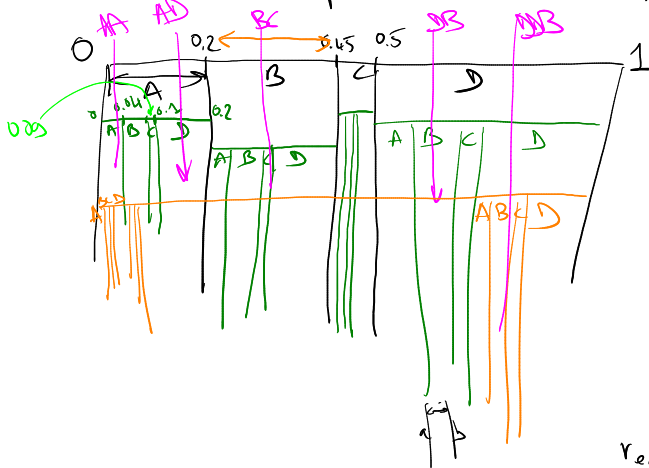
A: 7
G: 6



Arithmetic coding

$A = \{A, B, C, \dots\}$ alphabet

↳ with predefined order: A, B, C, D ...



Length associated to symbol 'B' is its frequency

draw a random number uniformly in $[0, 1]$

↳ will have 25% of chances to be in the 'B' bin

↳ can pick randomly symbols with the frequency of the distribution

read: the probability that a string starts with "DB" is the width of the bin coding for "DB"

1) Encoding

given string: DBAC → yields an interval

$[a, b]$ → choose the binary representation of a real number that falls into that bin, i.e. $x \in [a, b]$

2) Generation

pick randomly (uniformly) a real number $\in [0, 1]$ and decode it (its binary representation)

0.011011001...
 $\frac{1}{4} + \frac{1}{8} + \frac{1}{32} + \frac{1}{64}$ $\in [a, b]$ stops when no ambiguity anymore (between bins)

→ when to stop decoding (arbitrary number in binary format)

↳ use end-of-file character → cost?

$$H(p)$$

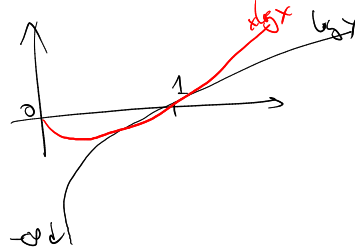
$$p' : \begin{cases} \dots \varepsilon \text{ for EOF} \\ 1-\varepsilon \text{ for standard symbols} \end{cases} \rightarrow \begin{matrix} p(A) \text{ for symbol } A \\ p(B) \dots \end{matrix} \quad \left. \begin{matrix} p'(A) = (1-\varepsilon)p(A) \\ p'(EOF) = \varepsilon \end{matrix} \right\}$$

$$H(p') = -\varepsilon \log \varepsilon - (1-\varepsilon) \log (1-\varepsilon) + (1-\varepsilon) H(p)$$

↳ choose ε : take the best for source data

continuous in ε

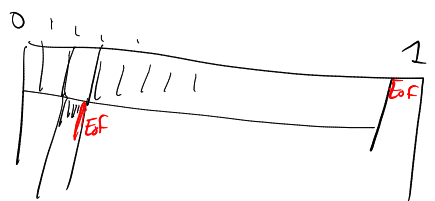
$$\approx H(p)$$



Ex for EOF (End of File)

Alphabet = { Latin characters + ' ' + punctuation marks }

"Hello world!" → "Hello world! EOF"



Huffman coding: is an approximation of arithmetic coding with the constraint that bin sizes should be powers of 2 (i.e. 2^{-n})

VI Additional remarks

Law over integers

code 1: encoding length for number n : $\log n + 2 \log \log n$

$$p(n) \propto 2^{-\log n - 2 \log \log n} = \frac{1}{n (\log n)^2}$$

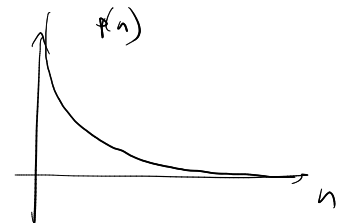
↳ the sum $\sum_n p(n)$ does converge → can be normalized as a probability distrib^o
↳ $\neq 1$

code 2: encoding length: $\log n + \log \log n + \log \log \log n + \dots$

$$\hookrightarrow p(n) \propto \frac{1}{n \times \log n \times \log \log n \times \dots}$$

↳ $\sum_n p(n)$ converge? → Kraft says yes

Limit case

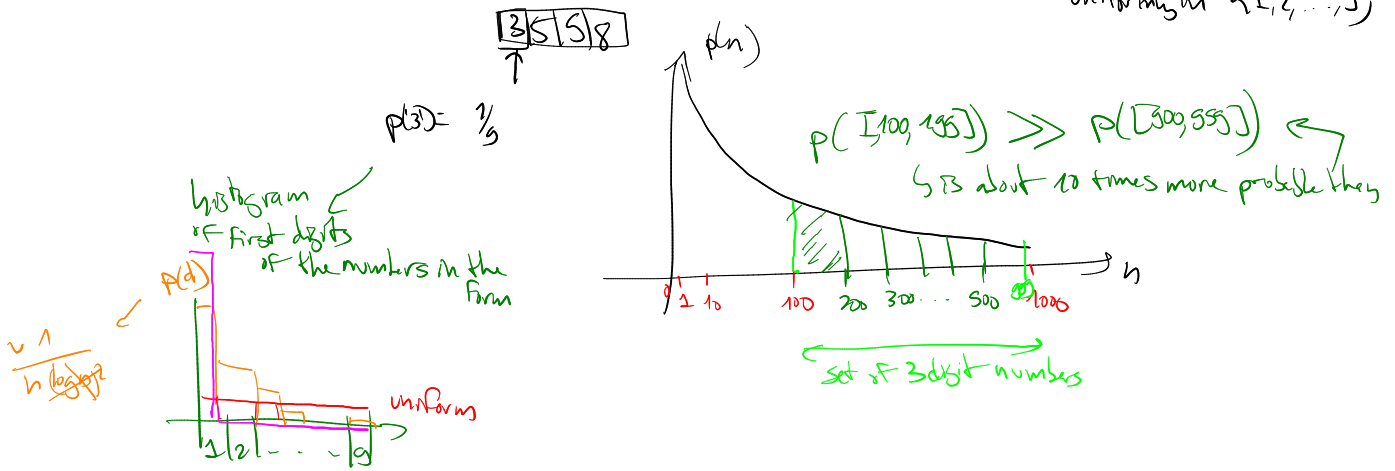


code? s.t. $p(n) \propto \frac{1}{n}$ or $\frac{1}{n \log n}$? → not possible because $\sum_n p(n) \rightarrow \infty$

Example: tax cheaters → fill the forms with fake numbers

↳ don't know information theory ...

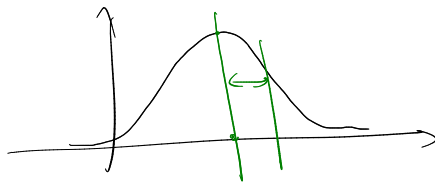
they pick the first digit randomly uniformly in $\{1, 2, \dots, 9\}$



• About generation : notion of typicality

data (made with n symbols) is typical of a model if its average coding length / symbol is close to the entropy of the model

encoding rate for that data



$$\mathbb{E}[(x-\mu)^2] = \sigma^2 \text{ standard deviation}$$

• Measure model efficiency or error : n bits

↳ Fit a model to a ML task → encoder

P → $-\log p$ unit: bits

VII Ex: text compression/prediction/generation with Markov Chains

→ of practical session

↳ IID model

↳ Markov chains with various orders

- compression = Hutter prize → prize hutter1.net

thousands € if able to compress file by 1% more
↳ Wikipedia

Conclusion

- good model for data : good for every task: prediction, generation, compression

- measure model suitability n bits → see generative models as compressors

↳ baseline: gzip (Lempel-Ziv)

- $-\log p$ is reachable with non-integer number of bits, thanks to arithmetic coding
not $\lceil -\log p \rceil$