# Chapter 3: Kolmogorov complexity

## Intro

Sequence completion : how?
       ''    probability? preference?
       ''    model?

Find a model for the following sequences!

-   010101010101010101 ....   :   repeat '01'
-   0110101000001001110011001100111111 ....    $\sqrt{2} - 1$ in binary
-   110111100111010110011101011011010101  :  too many '1' for uniform law

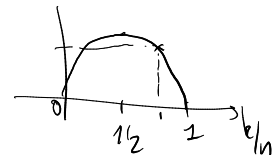$-\frac{1}{2}\lg\frac{1}{2} - \frac{1}{2}\lg\frac{1}{2} = -\lg\frac{1}{2}$
$= +\lg 2$
$= 1$

$0 \quad 1$
$\{0.5 ; 0.5\}$

English text + same one in German
  → compress each other independently : if short text
  → if very long (Full library):
       - encode how to translate English to German
       - encode the English text only

$k$ = count number of 1s

↳ length (encoding): $n \times H(p) = n$

encode: $\log n + n\,H\!\left(\frac{k}{n}\right)$

or
order all sequences
of length $n$ with $k$ 1s
then encode the index
   of this particular sequence

$\quad 0 \qquad 1$
$p: \frac{n-k}{n} \quad \frac{k}{n}$
$H\!\left(\frac{k}{n}\right) := H(p)$

$H\!\left(\frac{k}{n}\right)$



## I   Kolmogorov complexity

### Definition    [Cover & Thomas, p 463]

- $K(\text{string } s)$ = description cost of $s$ = the length of the shortest program that can produce it

  ↳ unit: bits

  ↳ ex:   $s = 00000000...00$   (1000 of them) : repeat 1000 '0' → 15
                                      $0^{1000}$ → 6

       $s = 3.14159...$    (1000 first digits of $\pi$)

- Associated probability:   $2^{-K(s)}$   → distribution over strings

  ↳ sum up to 1?   <u>no</u>

       $\sum_s 2^{-K(s)} \to \Omega$ : Chaitin's constant → provably non-computable
                                                ↳ halting problem

### Universality    [Cover & Thomas p. 427]

- more or less Turing-machine-invariant
  → universal Turing machine
  → Church's thesis: all computational models are equivalent (if sufficiently complex)

  $K_{\text{Computer 1}}(s) \le K_{\text{Computer 2}}(s) + C_{1,2} \quad \forall s$

       of size $C_{1,2}$
  ↳ by running a simulator of Computer 2 on computer 1
  and executing the program associated to $K_{\text{Computer 2}}(s)$

  ⇒ we'll always have + constant
              in all formulas bounding $K(s)$

- not dependent on encoding:
  - given 2 possible binary encodings $f$ & $g$ for data $x$ (which is not binary)

$$K(g(s)) \leq K(f(x)) + K(g \circ f^{-1}) \quad (+ c^{st})$$

<span style="color:magenta">translation from an encoding to the other one</span>

- these constants are small : $< 1MB$

    compared to possible big data size $(GB, TB)$

↳ Kolmogorov complexity really makes sense

## Extensions (defined up to a multiplicative factor, $e^{+c^{st}}$)

$P_1$ - probability associated with Kolmogorov complexity : $p(s) = 2^{-K(s)}$

$P_2$ - other version : $\displaystyle\sum_{\substack{\text{all programs } f \\ \text{that produce } s}} 2^{-\text{length}(f)}$

$P_3$ - with probabilistic programs : $\displaystyle\sum_{\substack{\text{all random} \\ \text{programs } f}} 2^{-\text{length}(f)} \times P(f \text{ outputs } s)$

<span style="color:red">← probability</span>

$P_4$ - with distributions : $\displaystyle\sum_{\substack{\text{all proba} \\ \text{distributions } \mu}} 2^{-\text{length to describe}(\mu)} \times \mu(s)$

**Proposition** : these 4 definitions are equivalent, i.e. :

$\forall$ definitions $i, j, \exists c \in \mathbb{R}, \quad p_i \leq c \, p_j$

[Zvonkin & Levin, 1970]

<span style="color:red">$\exists c_1, c_2 > 0$</span>
<span style="color:red">$\forall s : c_1 \, P_1(s) \leq P_2(s) \leq c_2 \, P_1(s)$</span>

↳ named <u>Solomonoff</u> universal prior (for prediction)

<u>Relative complexity</u>: $K(s|z)$ when $z$ is already available → no need to describe it

↳ similar to mutual information, conditional entropy, etc.

## II Bounds

### Easy upper-bound

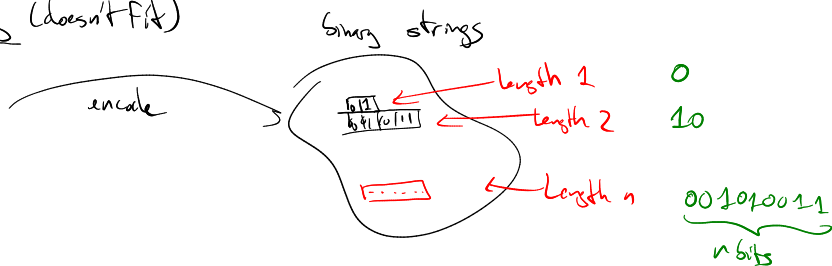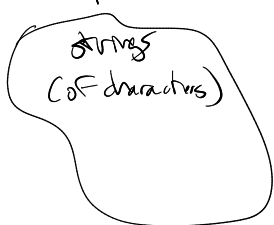· $K(s) \leq \text{length}(s) + 2\log \text{length}(s) + c$

<span style="color:green">or $+ \log + \log\log + \log\log\log + \ldots + c$</span>
<span style="color:green">(cf previous course)</span>

Program : "print the following chain, of length ___ : ___"

<span style="color:red">↑ length(s)   ↑ s</span>

· $K(s) \leq |s| + K(\boxed{|s|}) + c$ <span style="color:red">more general</span>

<span style="color:red">↓ ∈ ℕ</span>

· <u>Can't compress everything</u> (doesn't fit)

strings (of characters) → encode → Binary strings

<span style="color:red">length 1</span> <span style="color:green">0</span>
<span style="color:red">length 2</span> <span style="color:green">10</span>
<span style="color:red">length n</span> <span style="color:green">001010011</span>
<span style="color:green">n bits</span>

number of strings $s$
s.t. $K(s) \leq n$
is $1 + 2 + 2^2 + 2^3 + 2^4 + \ldots + 2^n < 2^{n+1}$
⟹ not many strings are simple

$\hookrightarrow$ strings $s$ such that $K(s \mid |s|) \geq |s|$ are named "algorithmically random" by <u>Kolmogorov</u>
(because no regularity to exploit)

$\hookrightarrow$ infinite strings $s = (x_1, x_2, \ldots)$ such that $\lim\limits_{n\to\infty} \dfrac{K(x_1 x_2 \ldots x_n (n))}{n} = 1$ are called <u>incompressible</u>

. $K(s) \leq |zip(s)| + |unzip\ program|$     "Clustering by compression" [Cilibrasi Vitangi]

$\hookrightarrow$ distance based on zip used to cluster files (text, MIDI ...) $\longrightarrow$ it worked!

$\hookrightarrow d(x,y) = \dfrac{\max(K(x|y), K(y|x))}{\max(K(x), K(y))}$     using zip as a proxy for $K$

. If $x$ is in $\Sigma$ (finite set):
$$K(x) \leq K(\Sigma) + \lceil \log |\Sigma| \rceil + c$$

<span style="color:red">$\uparrow$ say which element of the set $\Sigma$ $x$ is</span>

$\hookrightarrow$ <u>generalization</u>: given a set $X$ (not necessarily finite) and a probability measure on it,
$$K(x) \leq K(\mu) - \log \mu(x) + c$$

$\hookrightarrow$ in machine learning: a model $\mu$ is good if this quantity is small!

. <u>Theorem</u>:  <u>Kolmogorov complexity is non computable!</u>

- cf Gödel, Turing, halt problem $\rightarrow$ indecidability of the output of some programs
$\hookrightarrow$ whether they halt
$\hookrightarrow$ and if yes, what they output

- Cannot prove that $K(x) > 1 MB$, whatever $x$ is

"Berry paradox":
" The smallest number that cannot be described in less than 13 words. "     $\rightarrow$ description of a number using 12 words

"The first $x$ found that cannot be described in less than $L$ bits"

- Paradoxal proof in 2 steps, by chaitin:  <u>Incompleteness theorem</u>, 1971

- <u>step 1</u>: <u>proposition</u>:  $\exists L$ s.t. it is not possible to prove the statement $K(x) > L$ for any $x$

proof: $L = 1 MB$.
- write a program that goes through all possible proofs
and stop when finding a proof of $K(x) > L$ for some $x$
and print that $x$
- this program has length $< L$
- if it stops and prints an $x$ ... $x$ can be described by a program of length $< L$!
$\Rightarrow$ contradiction
$\Rightarrow$ this program never stops
$\Rightarrow$ there doesn't exist any proof of the form "$K(x) > 1 MB$" for any $x$

- <u>step 2</u>: <u>theorem</u>:  <u>Kolmogorov complexity is not computable</u>

proof: . consider all integers $\in [1, 2^{L+1}]$
. there are at most $2^{L+1} - 1$ programs of length $\leq L$
. $2^{L+1} > 2^{L+1} - 1 \Rightarrow \exists n_0 \in [1, 2^{L+1}]$ s.t. $K(n_0) > L$
. if there exists a program able to compute $K(n)$ for any $n$, by computing $K(n_0)$ we would prove $K(n_0) > L \Rightarrow$ not possible $\Rightarrow$ such a program does not exist.

$\Rightarrow$ not possible to have lower bounds on $K$ (sufficiently complex string)
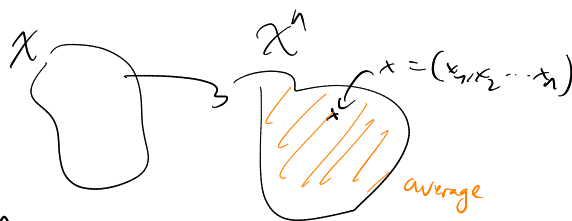↳ provided it's over 1 MB

## Entropic bound   [Cover & Thomas, p 473]

→ cannot prove a lower bound of $K$ for a given string $s$
   **but**
   can prove lower bounds **on average**

→ consider a proba distribution $\mu$ over a set $X$



$$H(X) \leq \underset{x \sim \mu^n}{\mathbb{E}}\left[\frac{1}{n}K(x)\right] \leq H(x) + (|X|+1)\frac{\log n}{n} + \varepsilon/n$$

(under the middle term) average $K$ complexity

(under the right term) $\xrightarrow[n \to \infty]{} 0$

$$\Rightarrow \boxed{\lim_{n \to \infty} \underset{x \sim \mu^n}{\mathbb{E}}\left[\frac{1}{n}K(x)\right] = H(X)}$$

cannot do better than entropy (on average)

## Proof:

**lower bound:**
- $K(s)$ is the length of an encoding $\Rightarrow$ Gibbs inequality $\Rightarrow$ cannot be better than entropy ($KL \geq 0$)

**upper bound:** by explicit construction → write a program outputting $s$ and check its length ($\geq K(s)$)

- encode $n$ : costs $2 \log n$
- string $s = x = (x_1, x_2 \cdots x_n)$ ← symbols from an alphabet $\mathcal{A} = \{'a', 'b', \ldots 'z'\}$
  ↳ count the # of appearance of each symbol in $s$ → $n_a \quad n_b \quad n_z$

- encode them:
  cost: $|A| \log n$ , actually $(|A|-1)$
  ↳ $|X|$

- now: draw the list of all possible sequences of $n$ characters with exactly $n_a$ 'a', $n_b$ 'b', ...
  $\underbrace{\text{of this list: size} \leq 2^{nH(\text{Bernoulli}(n_a/n, n_b/n, \ldots))}}$

  ↳ $C_k^n = \binom{n}{k} \leq 2^{nH(\ldots)}$ using Stirling formula

  ↳ $\sum_k \binom{n}{k} p^k (1-p)^{n-k} = 1$ with $p = k/n$
  $\underbrace{\quad}_{\leq 1}$

- encode the index of $s$ within this list
  ↳ $\leq n H(x)$   using Jensen

$\Rightarrow$ total encoding cost for $s$ = what is written in the formula

## III   Minimum Length Description principle (MDL)

**Def°** (General criterion for model selection
Given a set $X$ and a probability distribut° $\mu$ on it,
for any $x$, $\quad K(x) \leq K(\mu) - \log \mu(x) \quad (+c)$

(under $K(\mu)$) Complexity of the model

(under $-\log\mu(x)$) Likelihood: how well the model fits the data



Encoding cost w.r.t $\mu$:
$-\log p(x)$
$-\log \mu(x)$

= natural trade-off between model complexity and accuracy

Occam's razor
→ pick the model $\mu$ with the lowest $K(\mu) - \log \mu(x)$

<u>Unsupervised learning task</u>
dataset $\mathcal{D}$,  goal: generate new points according to $\mathcal{D}$
$\quad\hookrightarrow$ generator $G \to$ defines a probability $\boxed{P_G}$ over $\mathcal{X}$
$\qquad$ model $\qquad\qquad\qquad\qquad\qquad\qquad\qquad \mu$

$$K(\mathcal{D}) \leq K(G) - \underbrace{\sum_i \log P_G(x_i)}$$

$$-\log \prod_i P_G(x_i) = -\log P_G(\mathcal{D})$$

<u>Supervised learning task</u>
training dataset $\mathcal{D} = \langle (x_i, y_i) \rangle$
$\quad\hookrightarrow$ space $\mathcal{X} \times \mathcal{Y} =: \mathcal{Z}$
Goal: a function $\mathcal{X} \to \mathcal{Y}$
$$\qquad\qquad x \longmapsto y$$

deterministic view

Graph of the function
prior measure defined on $\mathcal{Z}$
$\quad\hookrightarrow P_M((x, y))$ ?$\qquad P_M(y|x) P_M(x)$
$\qquad\qquad\qquad\qquad\qquad$ function

non-deterministic view

$$K(\mathcal{D}) \leq K(P_M) - \log P_M(\mathcal{D})$$
$$-\log \prod_{i \in \mathcal{D}} P_M(x_i, y_i)$$
$$= \sum_i -\log P_M(x_i, y_i)$$

Examples
- Dirac peak. $\mu = \delta_{\mathcal{D}}$
$\quad\hookrightarrow \underbrace{K(\mu)}_{?} \underbrace{-\log \mu(\mathcal{D})}_{0}$
$\qquad\quad K(\mathcal{D})$

Space of all possible datasets
$\mathcal{D}$
$\mu$

overfit: high model complexity,
high likelihood

- Gaussian distribution $\mathcal{N}(m, \sigma)$, fit to a cloud of points $\mathcal{D} = \langle x_i \rangle$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \Sigma$

bigger $\sigma$
$\langle m \rangle \Sigma$
$\mathcal{X}$

$$K(\mu) - \log \mu(\mathcal{D})$$
$$-\sum_i \log P_\mu(x_i)$$
$$\frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(x_i - m)^2}{2\sigma^2}}$$
(in dimension 1)

$K(m)$
$K(\sigma)$

encode 2 real
values to encode
$\downarrow$
up to which precision?
$m = 1,205723301\ldots$

$$\sum_i + \frac{1}{2}\log(2\pi\sigma) + \frac{(x_i - m)^2}{2\sigma^2}$$

data term: $\sum_i (x_i - m)^2$
regularizer $\quad +$
factor?
$\quad\hookrightarrow$ Given by MDL
$\quad$ (no need to search for it)
$\quad\hookrightarrow$ MDL is a general approach
$\quad$ to define ML problems
$\quad$ i.e. to set up the training
$\qquad\qquad\qquad\qquad$ criterion

trade off between the precision
of $m$ encoded
& the accuracy of the model $\mathcal{N}(m, \sigma)$

<u>Model selection</u>
given $N$ models $\mu_k$, select model $\mu_{k^*}$ that minimizes $K(\mu) - \log \mu(\mathcal{D})$

## Instanciations

. AIC, BIC : approximations of $K$ (model)

$\hookrightarrow$ __AIC__: Akaike Information Criterion [1973]

$K(\mu) :=$ number of parameters of the model $\mu$

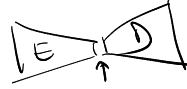$\hookrightarrow$ __BIC__: Bayesian Information Criterion [Schwartz 1978]

$K(\mu) := \frac{1}{2}$ number of parameters $\times \log($number of observations$)$

$\searrow$ dataset size

## Restricted Families of programs

- $K(s) := |zip(s)|$

- auto-encoders : the middle layer
    = compressed data



- programs on Turing machines
    $\longrightarrow$ restrict to finite automata $\implies$ Kolmogorov complexity is computable

## IV Conclusion

_ MDL = very general principle, to formulate any Machine Learning problem