

Complexity and Information Theory

Pietro Marcatti

First Semester 2022/2023

1 Introduction

The information communicated with a sentence, or any other way, is always dependant on the context. Same goes for the quality for the information. If an event has a low probability of occurring in a given context the information it provides is high, inversely if the probability is high the information is little. We could try to describe information with a simple inverse model:

$$Information = \frac{1}{p(E)}, \quad E=\text{event}$$

$$p(E) \rightsquigarrow 0 \quad Information \rightsquigarrow \infty$$

$$p(E) \rightsquigarrow 1 \quad Information \rightsquigarrow 1$$

Instead if we were to adopt a logarithmic model:

$$Information = \log \frac{1}{p(E)} = -\log p(E), \quad E=\text{event}$$

$$p(E) \rightsquigarrow 0 \quad Information \rightsquigarrow \infty$$

$$p(E) \rightsquigarrow 1 \quad Information \rightsquigarrow 0$$

Let's consider an alphabet as a group of possible events for which we have a probability distribution:

$$A = a_1, a_2, a_3, \dots, a_k$$

$$P = p_1, p_2, p_3, \dots, p_k$$

meaning that p_i is the probability to observe the event a_i . We can then calculate the average quantity of information as follows:

$$\sum_{i=1}^k p_i \cdot (-\log p_i) = \underset{\text{ShannonEntropy}}{\mathcal{H}(P)} \quad (1)$$

1.1 Coin Flip Example

FAIR COIN

$$A = \{H, T\}$$

$$P(H) = \frac{1}{2}, \quad P(T) = \frac{1}{2}$$

$$\mathcal{H}(P) = \frac{1}{2} \cdot (-\log \frac{1}{2}) + \frac{1}{2} \cdot (-\log \frac{1}{2}) = 1$$

Every time we flip the fair coin we expect to gain a bit of information. UNFAIR COIN

$$A = \{H, T\}$$

$$P(H) = \frac{9}{10}, \quad P(T) = \frac{1}{10}$$

$$\mathcal{H}(P) = \frac{9}{10} \cdot (-\log \frac{9}{10}) + \frac{1}{10} \cdot (-\log \frac{1}{10}) \rightsquigarrow 0.32$$

1.2 Properties of Entropy

- $\mathcal{H}(P) \geq 0$
- $\mathcal{H}(P) = 0$ if there is an event with probability 1
- \mathcal{H} is continuous with respect to P
- If an event gets split the entropy should be additive
- $\mathcal{H}(p_1, p_2, \dots, p_k) \leq \log k = \mathcal{H}(\frac{1}{k}, \frac{1}{k}, \dots, \frac{1}{k})$

To prove the last property we can first show that for a equi-distribution we get $\mathcal{H} = \log k$. To prove that this is also the maximum we use Jensen disequality:

$$\forall f, f''(x) < 0 \quad \text{in } [a, b], (a, b) \in \mathbb{R}^2$$

$$f\left(\sum_{i=1}^k \lambda_i \cdot x_i\right) \geq \sum_{i=1}^k \lambda_i \cdot f(x_i), \quad \sum \lambda_i = 1$$

2 Data Compression

Before we can talk about data compression it is necessary to explore the meaning and the functioning of encodings. Given two alphabets $A = a_1, a_2, \dots, a_k$ and $B = b_1, b_2, \dots, b_k$ an encoding is function:

$$\varphi : A^* \longrightarrow B^* \quad (2)$$

For a function to be a suitable encoding it must be at least an injective function (one-way function). In information theory and in data encoding in particular we refer to injective encoding as **uniquely decodable** encodings. Prefix codes are uniquely decodable codes that can be decoded without delay and are written in the form:

$$A = \{a_1, a_2, \dots, a_k\} \xrightarrow{\varphi} B = \{b_1, b_2, \dots, b_D\} \quad (3)$$

where $\varphi_i = |\varphi(a_i)|$

2.1 Kraft-McMillen 1958 Inverse Theorem

If φ is a uniquely decodable code, then:

$$\sum_{i=1}^k D^{-l_i} \leq 1 \quad (4)$$

D is the cardinality of the output alphabet, l_i is the length of the encoding for a_i . We can quickly see that, there can't be an uniquely decodable code that uses three encodings of length 1 to map over the binary alphabet:

$$A = \{a, b, c\}, \quad B = \{0, 1\}, \quad l_a = l_b = l_c = 1 \longrightarrow \frac{1}{2} + \frac{1}{2} + \frac{1}{2} > 1 \quad (5)$$

Instead, if we were to insist on using a length 10 encoding for everyone of the input characters:

$$A = \{a, b, c\}, \quad B = \{0, 1\}, \quad l_a = l_b = l_c = 10 \longrightarrow \frac{1}{2^{10}} + \frac{1}{2^{10}} + \frac{1}{2^{10}} < 1 \quad (6)$$

2.1.1 Proof 1

Assume the code is uniquely decodable and it is a prefix code, we can sort the input alphabet so that the length of the encoding satisfies the following:

$$a_1 \leq a_2 \leq \dots \leq a_k$$

A character with an encoding of length i generates a sub-tree rooted in a_i with height $l - l_i$. This sub-tree contains $D^{(l - l_i)}$ nodes that cannot be used in the encoding because they would share an encoded prefix.

Insert diagram

In the original three the following number of leaves

$$D^l - \sum_{i=1}^{k-1} D^{l-l_i}$$

The $k - 1$ term that bounds the sum is needed because we need to make sure that one last leaf is available to assign it to the last character a_k .

$$\begin{array}{lcl} D^l - \sum_{i=1}^{k-1} D^{l-l_i} \geq 1 & & \\ D^l(1 - \sum_{i=1}^{k-1} D^{-l_i}) \geq 1 & \left. \begin{array}{l} \text{Group by } D^l \\ \text{Divide by } D^l \end{array} \right\} & \\ 1 - \sum_{i=1}^{k-1} D^{-l_i} \geq D^{-l} & & \\ 1 \geq \sum_{i=1}^k D^{-l_i} & \left. \begin{array}{l} \\ D^{-l} \text{ is the last term of the sum} \end{array} \right\} & \end{array}$$

2.1.2 Proof 2 - General case

Let's define $N(n, h)$ the number of strings of alphabet A^n having encoding length h . It must be true that $N(n, h) \leq D^h$ because φ is uniquely decodable.

$$D^{-l_1} + D^{-l_2} + \dots + D^{-l_k} \leq 1 \quad (7)$$

$\forall n, n \in \mathbb{N}$ let's consider the object $(D^{-l_1} + D^{-l_2} + \dots + D^{-l_k})^n$. If written as a product it takes the following form:

$$D^{-l_1 \cdot n} + D^{-l_1 \cdot (n-1) - l_2} + \dots + D^{-l_k \cdot n}$$

As n tends towards infinity the exponent will behave differently depending on the value of the base. If the base is < 1 it will be bound between 0 and 1, and in any case < 1 . If the base is > 1 it will grow towards infinity faster than any polynomial function. Lastly, if the base is $= 1$ it will remain constant. To prove the theorem we will demonstrate that the object is dominated by a linear function, meaning that it cannot be in the second case. All of the members' exponents follow this chain of disequalities: $l_1 \cdot n \leq \text{exp} \leq l_k \cdot n$. That is, because we sorted the character by encoding length, 1 being the shortest while k the longest.

$$\begin{aligned} N(n, 1)D^{-1} + \dots + N(n, l_1 \cdot n)D^{l_1 \cdot n} + \dots + N(n, l_k \cdot n)D^{-l_k \cdot n} \\ \text{prob.0} \\ \leq D^1 D^{-1} + \dots + D^{l_k} D^{-l_k} \leq l_k \cdot n \end{aligned}$$

Since the object is dominated by $l_k \cdot n$ which is linear, it must mean that it also grows at most linearly, meaning that the base of the exponent must be ≤ 1 .

2.2 Kraft-McMillen Direct Theorem

If given l_1, l_2, \dots, l_k and D such that $\sum_{i=1}^k D^{-l_i} \leq 1$ there must exist a prefix code φ having l_1, l_2, \dots, l_k as lengths of the encoding.

2.2.1 Proof

Proof can be given by construction by producing the complete D-ary tree

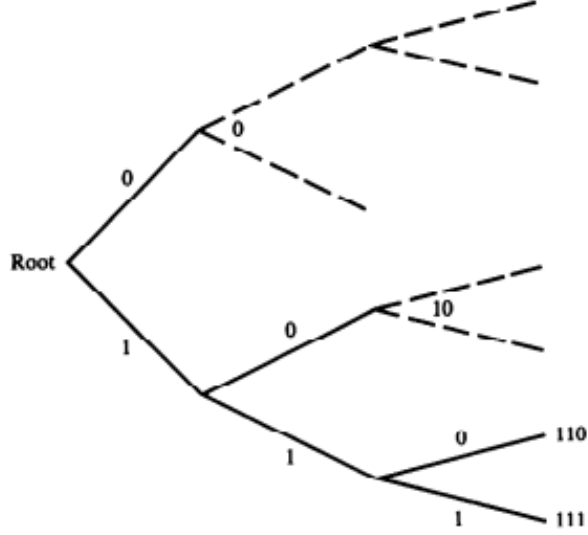


Figure 5.2. Code tree for the Kraft inequality.

Observation: Prefix codes compress as good as any other uniquely decodable code.

2.2.2 Average expected code length - Definition

Given the input alphabet A , P the probability and B the output we define the Expected Length as:

$$EL(\varphi) = \sum_{i=1}^k p_i \cdot l_i = \sum_{i=1}^k p_i |\varphi(a_i)| \quad (8)$$

Our goal is to find prefix codes that minimize EL .

2.3 1st Shannon Theorem 1948

If φ is uniquely decodable code then $EL(\varphi) \geq \mathcal{H}_D(P) = \sum_{i=1}^k p_i \cdot \log_D p_i$. The proof relies on Kraft-McMillen theorem and this simple logarithmic property:

$$\log_e x \leq x - 1, \quad -\log_e x \geq 1 - x$$

2.3.1 Proof

$$\begin{aligned}
EL(\varphi) - \mathcal{H}_D(P) &= \sum_{i=1}^k p_i \cdot l_i + \sum_{i=1}^k p_i \cdot \log_D p_i && \left. \vphantom{\sum_{i=1}^k p_i \cdot l_i} \right\} \text{Group by } p_i \text{ and force } l_i \text{ as } \log_D \\
&= \sum_{i=1}^k p_i \cdot \log_D(D^{l_i} p_i) && \left. \vphantom{\sum_{i=1}^k p_i \cdot \log_D(D^{l_i} p_i)} \right\} \text{Change to base } e, \text{ take out the common term} \\
&= \frac{1}{\log_e D} \cdot \sum_{i=1}^k p_i \cdot \log_e(D^{l_i} p_i) && \left. \vphantom{\sum_{i=1}^k p_i \cdot \log_e(D^{l_i} p_i)} \right\} \text{Group by } -1 \\
&= \frac{-1}{\log_e D} \cdot \sum_{i=1}^k p_i \cdot \log_e \frac{1}{D^{l_i} p_i} && \left. \vphantom{\sum_{i=1}^k p_i \cdot \log_e \frac{1}{D^{l_i} p_i}} \right\} \text{Using the above property} \\
&\geq \frac{-1}{\log_e D} \cdot \sum_{i=1}^k p_i \cdot \left(\frac{1}{D^{l_i} p_i} - 1 \right) \\
&= \frac{-1}{\log_e D} \cdot \left(\sum_{i=1}^k \frac{1}{D^{l_i}} - \sum_{i=1}^k p_i \right) && \left. \vphantom{\sum_{i=1}^k \frac{1}{D^{l_i}}} \right\} \text{Group by } -1 \\
&= \frac{1}{\log_e D} \left(1 - \sum_{i=1}^k \frac{1}{D^{l_i}} \right) \underset{KMM \rightarrow \geq 0}{\geq} 0
\end{aligned}$$

An optimum code achieves the equality between the expected length and the entropy of the distribution.

Shannon Code

Let's recall the definition of expected length and entropy over the probability P and alphabet with cardinality D :

$$EL(\varphi) = \sum_{i=1}^k p_i \cdot l_i$$

$$\mathcal{H}_D(P) = \sum_{i=1}^k p_i \cdot (-\log_D p_i)$$

Observation: $\log_D \frac{1}{p_i} \in \mathbb{R}$, $l_i \in \mathbb{N} \longrightarrow l_i = \lceil \log_D \frac{1}{p_i} \rceil$

If $\sum_{i=1}^k D^{-l_i} \leq 1$ then we can use the Direct Theorem to prove φ , with encoding lengths l_1, l_2, \dots, l_k , exists. The encoding is obtained with a greedy strategy on the D -ary tree:

$$\sum_{i=1}^k D^{-l_i} = \sum_{i=1}^k D^{-\lceil \log_D \frac{1}{p_i} \rceil} \leq 1$$

$$\lceil \log_D \frac{1}{p_i} \rceil = \log_D \frac{1}{p_i} + \beta_i \quad 0 \leq \beta_i < 1$$

$$\sum_{i=1}^k D^{\log_D p_i - \beta_i} = \sum_{i=1}^k D^{\log_D p_i} \cdot \frac{1}{D^{\beta_i}} = \sum_{i=1}^k p_i \cdot \frac{1}{D^{\beta_i}} \leq 1$$

Example: $A = \{a, b\}$ $B = \{0, 1\}$ $P = \{1 - \frac{1}{32}, \frac{1}{32}\}$ $l_b = \log_2 2^5 = 5$ $\lceil \log_2 1 - \frac{1}{32} \rceil = 1 = l_a$ But this is not efficient because we are assigning unnecessarily long codes to characters with low probability.

Sub-Optimality of Shannon Codes

We can demonstrate that Shannon Codes are suboptimal:

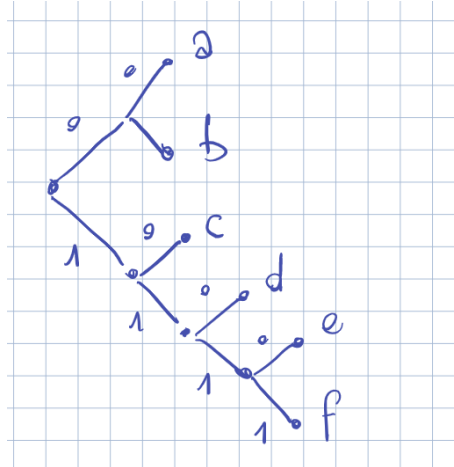
$$\begin{aligned}
 EL(S) &= \sum_{i=1}^k p_i (-\log p_i + \beta_i) \\
 &= \sum_{i=1}^k p_i (-\log p_i) + \sum_{i=1}^k p_i \beta_i \\
 &= \mathcal{H}_D(P) + \sum_{i=1}^k p_i \beta_i \leq \mathcal{H}_D(P) + 1 \quad \sum_{i=1}^k p_i \cdot \beta_i \leq 1 \\
 \mathcal{H}_D(P) &\leq EL(\varphi) \leq \mathcal{H}_D(P) + 1
 \end{aligned}$$

Shannon-Fano Code

Given the alphabet A with characters a_1, a_2, \dots, a_k mapping to the probabilities $p_1 \leq p_2 \leq \dots \leq p_k$. The aim is to balance the sum of the probabilities assigned to the branches of the root. To obtain balance we try to minimize this absolute value: $|\sum_{i=1}^h p_i - \sum_{i=h+1}^k p_i|$. We then repeat the calculation recursively. Example:

$$A = \{a, b, c, d, e, f\}$$

$$P = \left\{ \underbrace{\frac{40}{100}, \frac{18}{100}}_{\frac{58}{100}}, \underbrace{\frac{15}{100}, \frac{13}{100}, \frac{10}{100}, \frac{4}{100}}_{\frac{42}{100}} \right\}$$



Shannon-Fano is compromising on splitting the characters set in two of equal

(or best possible) probability because the problem of splitting a set in equally weighted subsets is NP-Complete.

Shannon on strings of lenght n over input alphabet

Example: $A = \{a, b\}$ $B = \{0, 1\}$ $P = \{\frac{3}{4}, \frac{1}{4}\}$ $\mathcal{H}_2(P) = 0.81$ $EL(\varphi) = 1.25$
Let's define the efficiency of a code as follows:

$$Eff. = \frac{\mathcal{H}_D(P)}{EL(\varphi)} \leq 1$$

Applying to our example we get:

$$Eff. = \frac{\mathcal{H}_2(P)}{EL(\varphi)} = 0.64$$

Let's try to calculate the efficiency now with an encoding for pairs of characters:

$$A' = \{aa, ab, ba, bb\} \quad P' = \left\{ \frac{9}{16}, \frac{3}{16}, \frac{3}{16}, \frac{1}{16} \right\} = P^2$$

$$\mathcal{H}_D(P') = 1.62 = \mathcal{H}_D(P^2) = 2 \cdot \mathcal{H}_D(P)$$

$$EL(\varphi) = 1.93 \quad Eff. = 0.83$$

If we pretend we want to send a message 1000 characters long, with the single character encoding we expect a message lenght of 1.25k characters, Instead with the pairs encoding we expect $\frac{1.93}{2} \rightsquigarrow 0.97k$ characters.

Property: given two events x, y independent $\mathcal{H}(x \wedge y) = \mathcal{H}(x) + \mathcal{H}(y)$ Proof:

$$\begin{aligned} \mathcal{H}(x \wedge y) &= - \sum_{i,j} p_i \cdot q_j \log(p_i \cdot q_j) = - \sum_{i,j} p_i \cdot q_j \log p_i - \sum_{i,j} p_i \cdot q_j \log q_j \\ &= - \sum_{j=1} q_j \cdot \sum_i p_i \log p_i - \sum_{i=1} p_i \cdot \sum_j q_j \log q_j = \mathcal{H}(x) + \mathcal{H}(y) \end{aligned}$$

From Shannon Theorem we then have $EL(\varphi_n) \geq \mathcal{H}_D(P^n) = n \cdot \mathcal{H}_D(P)$

For Shannon Codes (and all sub-optimal codes) we have:

$$n \cdot \mathcal{H}_D(P) \leq EL(\varphi_n) < n \cdot \mathcal{H}_D(P) + 1$$

If we want to reason in terms of a single input alphabet value we divide the disequality by n. We can see that for $n \rightsquigarrow \infty$ the expected lenght is squeezed between the terms of the disequality.

Huffman Codes

Given the input alphabet and the probability distribution $A = \{a_1, a_2, \dots, a_k\}$ $P = p_1 \leq p_2 \leq \dots \leq p_k$ We start by taking the two least probable characters and we put them at the leaves, creating a phantom character with probability the sum of the two. I keep repeating this process until we obtain a complete tree.

History

- 1949 Shannon-Fano code
- 1952 Huffman code
- 1970 Huffman becomes popular
- 1977 Lempel and Ziv create the LZ77 Code (Variable length - Block) Uses a dynamic dictionary for compression.
- 1978 They update it with a new version LZ78 with static dictionary Gets patented by Sperry Corporation → Unisys 1981
- 1983 Another variant LZW (W - Welch) is a simpler implementation of LZ78. In the same year they request a patent and publish an article in 1984 without mentioning its patent.
- 1985 the patent gets actually officialized
- Unix and Compuserve start using the code
- 24/12/1994 Unisys decides to claim the royalties "The GIF Tax"
- 2004 the patent expires
- Unix decides to avoid the taxation by "downgrading" to LZ77

LZ-77

The idea is to focus on a part of the message i've already codified and i refer to it as search buffer. The portion just in front of the search buffer is the look-ahead buffer. The length of both buffers is a parameter of the implementation. I then operate like this: I look for the longest prefix of the look-ahead buffer that also appears as a substring in the search buffer. In reality it's a little more complex as the substring from the search buffer can overflow in the look-ahead buffer. Let's identify as j the position occupied by the first character of the substring and i the position of the first character of the prefix. We then consider the triple (o, l, s) :

- o : distance between i and j
- l : length of the prefix-suffix
- s : the first mis-matching character

Example

input = babbababbabbabba Buffer size: 5

Encoding:

1st step S buffer = \emptyset , LA buffer = babba, triple = (0,0,b)

2nd step S buffer = b, LA buffer: abbab, triple (0,0,a)

3rd step S buffer: ba, LA-buffer = bbaba, triple= (2,1,b)

4th step S buffer: babb, LA-buffer = ababb, triple = (3,2,a)

5th step S buffer: bbaba la-buffer = bbabb, triple = (5,4,b)

6th step S buffer:bbabb la-buffer= abba, triple= (3,4,@)

Decoding:

1st step @@@@b

2nd step @@@@@ba

3rd step @@@@@babb

4th step @@@@@babbaba

5th step @@@@@babbababbabb

6th step @@@@@babbababbabbabba

Complexity

Kolmogorov

History

We've talked about Shannon who worked in the US, studied at MIT and Princeton, collaborated with Bell Labs. We've talked about Fano who started at the Politecnico di Torino and then moved to MIT. Then we talked about Huffman who worked at Ohio State University, then founded the University of California Santa Cruz.

Kolmogorov, unlike the other names we've seen so far, does not come from the west. In 1965 he was working at the University of Moscow

Idea

Consider a computational model and define the complexity of a string as the length of the shortest program (in that computational model) that can generate that string. We'll use Turing Machines: $m = (k, \Sigma, \delta, s)$

k : finite set of states

Σ : finite alphabet $\triangleright, \sqcup \in \Sigma$

$\delta : K \times \Sigma \longrightarrow (K \cup \{yes, no, halt\}) \times \Sigma \times \{\rightarrow, \leftarrow, -\}$

$s : s \in K$ start state

Church-Turing Thesis

All computational models are turing-equivalent.

Turing machines with k tapes and I/O $m = (K, \Sigma, \delta, s)$

- $\delta : \Sigma^k x K \rightarrow \Sigma^k x (K \cup \{yes, no, halt\}) x \{\rightarrow, \leftarrow, -\}^k$
- The input tape (1st tape) cannot be modified but we can move backwards as much as possible
- The output tap (last tape) can only go forwards

Universal Machine Theorem

$$\exists u, u \text{Universal Turing Machine} | u(bin(m), x) = m(x)$$

How does the universal turing machine works: its tape is something like this: start $[01 \dots 1010]$ x To be able to go on with the computation if we were in the middle of it, we need to know the configuration of the machine. The configuration is stored in one of the k tapes. We also need the state we are on and the position of the first tape. The tape and the position can be expressed as a triple (q, w, ς)

- q is the current character
- w is the string left of q
- φ is the string right of q

The Kolmogorov complexity of a string m is denoted by $K_u(x) = \min_{u(bin(m))=x} |bin(m)|$

Observation

Unfortunately this notion is not computable.

$$\nexists A \text{ Turing Machine} \quad | \quad A(x) = K_u(x)$$

Conditional Kolmogorov Complexity

The Conditional Kolmogorov Complexity of x given y is:

$$K_u(x|y) = \min_{u(bin(m),y)=x} |bin(m)|$$

This means that, obviously, if we provided some more information $K_u(x|y) \leq K_u(x)$. One of the most common information given is $|x|$.

Theorem 2.1 (Universality of Kolmogorov Complexity) *If \mathcal{U} is a universal computer (Turing Machine), for any other computer \mathcal{A} there exists a constant $c_{\mathcal{A}}$ such that*

$$K_{\mathcal{U}}(x) \leq K_{\mathcal{A}}(x) + c_{\mathcal{A}}$$

for all strings $x \in 0,1^*$, and the constant $c_{\mathcal{A}}$ does not depend on x .

Proof 2.4 Assume that we have a program $p_{\mathcal{A}}$ for computer \mathcal{A} to print x . Thus, $\mathcal{A}(p_{\mathcal{A}}) = x$. We can precede this program by a simulation program $s_{\mathcal{A}}$ which tells the computer \mathcal{U} how to simulate computer \mathcal{A} . Computer \mathcal{U} will then interpret the instructions in the program for \mathcal{A} , perform the corresponding calculations and print out x . The program for \mathcal{U} is $p = s_{\mathcal{A}}p_{\mathcal{A}}$ and its length is

$$l(p) = l(s_{\mathcal{A}}) + l(p_{\mathcal{A}}) = c_{\mathcal{A}} + l(p_{\mathcal{A}})$$

where $c_{\mathcal{A}}$ is the length of the simulation program. Hence, for all strings x

$$K_{\mathcal{A}}(x) = \min_{p:\mathcal{A}(p)=x} l(p) \leq \min_{p:\mathcal{U}(p)=x} l(p) + c_{\mathcal{A}} = K_{\mathcal{U}}(x) + c_{\mathcal{A}}$$

The crucial point is that the length of this simulation program is independent of the length of x , the string to be compressed. For sufficiently long x , the length of this simulation program can be neglected, and we can discuss Kolmogorov complexity without talking about the constants.

Theorem 2.2 (Upper bound for conditional kolmogorov complexity)

$$K(x \mid l(x)) \leq l(x) + c$$

Proof 2.5 Let's define a program for printing x as the program that says:

Print the following l -bit sequence: $x_1x_2\ldots, x_{l(x)}$

Note that no bits are required to describe l since l is given. The program is self-delimiting because $l(x)$ is provided and the end of the program is thus clearly defined. The length of this program is $l(x) + c$

An apparently stronger conclusion is found in the book "An introduction to Kolmogorov Complexity and its applications" by M. Li and P. Vitanyi, even though this relies on an important detail. In fact, in this case, the alphabet for the machine is $0, 1, \sqcup$. If we allow ourselves to make this assumption we can write the following.

Theorem 2.3 (Upper bound for kolmogorov complexity - Li-Vitanyi)

$$K(x) \leq l(x) + c$$

$$K(x) \leq K(x \mid l(x)) + 2 \log l(x) + c$$

Proof 2.6 IF the computer does not know $l(x)$ we must have some way of informing the computer when it has to come to the end of the string of bits that describes the sequence. We describe a simple but inefficient method that uses a sequence 01 as a "comma".

Suppose that $l(x) = n$. To describe $l(x)$, repeat every bit of the binary expansion of n twice; then end the description with a 01 so that the computer knows that it has come to the end of the description of n . For example, the number 5 (binary 101) will be described as 11001101. This description requires $2\lceil \log n \rceil + 2$ bits. Thus, the inclusion of the binary representation of $l(x)$ does not add more than $2 \log l(x) + c$ bits to the length of the program, and we have the bound in the theorem.

Theorem 2.4 (Upper bound for kolmogorov complexity - Cover)

Study reference for this part is taken from:

- Elements of Information Theory, Cover-Thomas, chapter 8
- An introduction to Kolmogorov Complexity and its applications, M. Li and P. Vitanyi

Invariance of Kolmogorov Complexity

$$\forall x \quad |K_U(x) - K_A(x)| \leq c, \quad c \text{ constant}$$

Theorem

$$K(x) \leq |x| + c, \text{ Li, Vitanyi}$$

Proof

U Printing Universal Turing Machine, we add a bit of information: if the first digit of the input tape is 0, then it simulates the rest of the tape as it is a binary encoding of a turing machine. Else if the first digit is 1 the machine will print in output the remaining significant digits (not blank).

Proposition

The number of strings having Kolmogorov complexity $< h$ is $< 2^h$

Proof:

$$2 \rightsquigarrow 1, 2^2 \rightsquigarrow 2 \dots$$

$$2^0 + 2^1 + \dots + 2^{h-1} = 2^h - 1$$

These are the machines that can have a binary encoding of lenght at most $h - 1$

Corollary

There is at least one string x such that $K(x) \leq |x|$

Proof: For a given h

at most $2^h - 1$ strings with $KC < h$

There are 2^h strings of lenght h . At least 1 string of lenght h has a $KC \geq h$

$$\forall x \quad K(x) \leq |x| + c$$

$$\exists x \quad K(x) \geq |x|$$

Kolmogorov Encoding

$$\varphi_K(x) = \text{bin}(m) \quad \text{U.D.}$$

$$\varphi_k : A^* \longrightarrow \{m | m \text{ Turing Machine}\}$$

$\varphi_k(x)$ is the shortest machine that produces $\text{bin}(x)$. This encoding is not computable but it's UD. The Kolmogorov complexity is then: $K(x) = |\varphi_k(x)|$.

$$EL_n(\varphi_k) = \sum_{x \in A^n} p(x) \frac{|\varphi_k(x)|}{K(x)}$$

$$EL_n(\varphi_k) \geq \mathcal{H}(P^n) = n\mathcal{H}(P)$$

where P is the probability distribution over A . The first partial result then is:

$$\frac{E_n(K(x))}{n} \geq \mathcal{H}(P)$$

To set a bound from the other side we reason as follows: any U.D. code φ can be seen as a set of machines for producing strings. $\{\text{Decoder } \varphi ; \varphi(x) \}$ is an algorithm that produces x as output.

$$\forall x \quad K(x) \leq |\text{Decoder}(\varphi)| + |\varphi(x)|$$

If applied to Shannon Code:

$$\begin{aligned} E_n(K(x)) &\leq |\text{Decoder}\varphi| + \sum_{x \in A^n} p(x)|\varphi(x)| \\ &\leq |\text{Decoder}\varphi| + EL_n(\varphi) \\ &\stackrel{\text{Shannon } C.}{\leq} |\text{Decoder}\varphi| + \mathcal{H}(P^n) + 1 \quad \text{Shannon Un-optimality} \\ &= |\text{Decoder}\varphi| + n\mathcal{H}(P) + 1 \\ &\leq K(P) + n\mathcal{H}(P) + 1 \\ \frac{E_n(K(x))}{n} &\leq \frac{K(P)}{n} + \mathcal{H}(P) + \frac{1}{n} \end{aligned}$$

Computational Complexity

We will focus on decision problems. Other classes of problems that we will not explore are: Functional problems and optimization problems.

Decision Problems : $P : I \longrightarrow \{yes, no\}$

Functional Problems :

Optimization Problems :

Having a fast (polynomial) solution for an optimization problems implies that i have a fast one also for the Functional version and the decision version of the problem. The same goes inversely with the decision being hard, complex (non polynomial).

The Computational model we will use is Turing Machines, read chapter 1 of the book.

Time Complexity

This is a decision problem. There's a problem P, a class of models of computation M. The goal is to find in M the fastest machine m that solves P. By fastest we mean that executes least instructions. Fix a notion of dimension of the input which is usually done in the definition of the model of computation.

COMPUTATIONAL CHURCH TURING THESIS:

All reasonable models of computation are polynomially related.

$$M_1 : \Theta(f(n)) \rightsquigarrow M_2 : \Theta(p(f(n))), \quad p \text{ polynomial}$$

We know already about turing machines but now we'll introduce Unlimited registry machines (URM):

Has an infinite set of registers (r_0, r_1, \dots, r_n) , in any of these registers there can be an arbitrarily long natural number.

- S(i) means that I sum 1 to the register i.
- Z(i) means that register i becomes 0.
- T(i,j) means that $r_j = r_i$.
- J(i,j,k) means that if the content of the $r_i = r_j$ jump to instruction k

Example:

P: compute $x+1$. For the turing machine this means receiving the binary digits of x and i want on output the result of $x+1$. The complexity is lineary with the number of digits.

With the URM the complexity is 1, i just need S(0). The problem is that we are hiding the size of the input so this model is not reasonable.

Let's make an addition to the URM model and add the instructions $P(i) = r_i = r_i * r_i$ we execute these instructions: $T(0,1)$, $J(1,2,6)$, $P(0)$, $S(2)$, $J(3,3,2)$ The output goes: $x, x^2, x^4, x^8, \dots, x^{2^x}$ with complexity $\Theta(n)$. With a turing machine only writing the digits of the output requires at least $\Theta : 2^x \log(x)$.

UNIFORM COMPLEXITY COST:

Each instruction of the models has complexity $\Theta(1)$. This is not reasonable when we are using models involving operations that make the involved integers grow too fast. An example of this behaviour is the URM with the product operation.

LOGARITHMIC COMPLEXITY COST:

Each instruction of the models has complexity that depends on the number of digits it manipulates.

It is fundamental to analyze the complexity of all the instructions in your computational model:

- $S(i): r_i = r_i + 1 \quad \Theta(\log(r_i))$
- $Z(i): r_i = 0 \quad \Theta(1)$
- $T(i,j): \Theta(\log(r_i))$
- $J(i,j,k): \Theta(\min(\log(r_i), \log(r_j)))$
- $P(i): \Theta((\log(r_i))^2)$

If we now analyze the cost of the models with the Logarithmic complexity cost we get that the URM and the Turing machine are related.

We define as $P = \{L | L \text{ can be decided in polynomial time on a deterministic Turing Machine}\}$.

P is invariant with respect to the model of computation (consequence of the extended Church thesis) We recall that a configuration for a Turing machine with k tapes is in the form $(q, u_1, w_1, \dots, u_k, w_k)$. The initial configuration is always: $(s, \triangleright, x, \triangleright, \sqcup, \dots, \triangleright, \sqcup)$ and a final configuration is always: $(H, u_1, w_1, \dots, u_k, w_k)$ where $H = \underbrace{\{yes, no\}}_{halt}$. For a machine that solves a decision problem L :

$$M(x) = yes \quad \text{iff} \quad x \in L \quad M(x) = no \quad \text{iff} \quad x \notin L$$

For a machine that does not terminate on input x we write $M(x) \uparrow$ and we say it diverges while the notation for a terminating TM is $M(x) \downarrow$.

The notion of computation: a computation step for a machine M is a binary relationship between two configurations:

$$(q, u_1, w_1, \dots, u_k, w_k) \longrightarrow (q', u'_1, w'_1, \dots, u'_k, w'_k)$$

(For reference Papadimitriou section 2.1)

Time Complexity

Definition: TIME COMPLEXITY

Let M be a k -TM and x be an input for M we say that M on input x takes time

t if:

$$(s, \triangleright, x, \triangleright, \varepsilon, \dots, \triangleright, \varepsilon) \longrightarrow (H, \dots)$$

with $H \in \text{yes}, \text{no}, \text{halt}$. We say that M operates in time at most $f(n)$ if:

$$\forall x \text{ with } |x| = m \quad M \text{ on } x \text{ takes time at most } f(|x|)$$

We will say that a language $L \in \text{TIME}(f(n))$ if $\exists M k\text{-TM}$ and M decides L and operates in time at most $f(n)$. $\text{TIME}(f(n)) = \{L, L', \dots\}$

Example: L is the language of all strings that are palindromes. $\Sigma = \{0, 1\} \cup \{\sqcup, \triangleright\}$ we get $\Theta(n^2)$

Example: M is a 2-TM we can solve the same problem in linear time. Is the quadratic loss shown here the worst possible loss between two TM machine models with different number of tapes?

THEO 2.1 PAPADIMITRIU:

if M is a k -TM that decides a language L in time $f(n)$, then there exists a 1-TM M' that decides L in time at most $\Theta(f(n)^2)$. fundamental hypothesis is that $f(n) \not\leq n$