
CODIFICACIÓN Y COMPRESIÓN

El concepto de entropía

Habitualmente se dice que la teoría de la información tiene su origen en un famoso artículo de C.E. Shannon de 1948 titulado “Una teoría matemática de la comunicación” [Sha48]. Entonces la informática no se había desarrollado (el famoso ENIAC entró en funcionamiento dos años antes y su “programación” se efectuaba conectando cables), sin embargo sí que tenía alguna importancia la transmisión de señales digitalizadas (el telégrafo muestra lo antigua que es la idea). Por ello la teoría introducida por Shannon es fundamentalmente digital y ha permanecido vigente a pesar de que el contexto de las comunicaciones ha cambiado totalmente. Es significativo de la influencia de [Sha48] que sea allí donde aparece por escrito por primera vez la palabra *bit* como abreviatura de *binary digit*, que Shannon atribuye a uno de sus colegas matemáticos pioneros de la informática.

Partimos de un conjunto finito $S = \{s_1, s_2, \dots, s_n\}$ que además es un espacio de probabilidad, con p_i la probabilidad de escoger s_i . La idea general es que los elementos de S sean cualesquiera fuentes de información (la notación S tomada de [Rom92] es de *source*) pero más adelante preferiremos considerar S como un conjunto de caracteres con los que componemos unidades mayores de información (que son un modelo para los ficheros digitales).

Shannon se preguntó en la primera parte de su artículo cómo medir la información contenida en S o dicho de otra forma más ilustrativa, la “incertidumbre” al extraer una muestra de S . Por ejemplo, si p_1 es casi 1 y el resto de los p_i son casi 0, apenas tenemos incertidumbre sobre qué elemento de S se va a elegir. El caso opuesto es que $p_1 = \dots = p_n = 1/n$ en el que no sabemos apostar por un elemento particular (véase en [Rén84] más sobre esta idea). Shannon procedió de una manera teórica y exigió a la función medidora de la información $H = H(p_1, \dots, p_n)$ que fuera continua, que $H(1/n, \dots, 1/n)$ fuera creciente en n (más oportunidades, más incertidumbre) y que cumpliera

$$(1) \quad H\left(\frac{1}{n}, \dots, \frac{1}{n}\right) = H\left(\frac{b_1}{n}, \frac{b_2}{n}, \dots, \frac{b_k}{n}\right) + \sum_{i=1}^k \frac{b_i}{n} H\left(\frac{1}{b_i}, \dots, \frac{1}{b_i}\right), \quad \forall b_i \in \mathbb{Z}^+ \quad \sum_{i=1}^k b_i = n.$$

Esta última propiedad simplemente significa que la cantidad de información no puede variar si subdividimos S en subconjuntos más pequeños de tamaño b_i [Rom92, §1.1] [Mac03, §2.5].

Shannon probó un pequeño resultado que asegura que las únicas funciones con estas propiedades son de la forma $H = K \sum_i p_i \log p_i$ con $K < 0$ y, como estaba interesado en el caso digital, tomó $K = -1/\log 2$. De esta forma, se define la *entropía* como

$$(2) \quad H(p_1, \dots, p_n) = - \sum_{i=1}^n p_i \log_2 p_i.$$

Por convenio, diremos que $p \log_2 p$ vale 0 para $p = 0$, pues es su límite cuando $p \rightarrow 0^+$. El nombre y la notación para la entropía vienen heredados de los trabajos de L. Boltzmann en el siglo XIX sobre un concepto análogo en mecánica estadística y termodinámica.

Si $n = 2^N$ y $p_1 = \dots = p_n = 1/N$, entonces la entropía es N , el número de bits que necesitamos para representar 2^N objetos. En general, si queremos representar un valor de k posibles necesitaremos aproximadamente $-\log_2(1/k)$ bits y (2) es un promedio de cantidades similares. Esto sugiere, como ya hemos indicado, que la entropía está relacionada con el número de bits necesarios para almacenar la información.

Para concretar esta idea, pensemos en que estamos interesados en comprimir textos y que en S incluimos sólo los caracteres y signos de puntuación del idioma que usamos. En esta situación es fácil tener una idea estadística de los valores de p_i , el propio Shannon aproximó con 4.14 la entropía de las letras en inglés. Para codificar un texto, concatenamos las codificaciones en bits de las letras. Definimos entonces *código (binario únicamente descifrable*, si se quiere ser más preciso) como una función inyectiva $C : S \rightarrow \{\text{cadenas de bits}\}$ que también es inyectiva cuando actúa sobre las cadenas de elementos de S

$$(3) \quad C : s_{i_1} s_{i_2} \dots s_{i_k} \mapsto C(s_{i_1}) C(s_{i_2}) \dots C(s_{i_k})$$

Por ejemplo, si partimos de $S = \{a, b, c\}$, entonces $C(a) = 0$, $C(b) = 1$, $C(c) = 10$ no es un código válido porque $C(ba) = C(c)$. Sin embargo sí lo es redefiniendo $C(b) = 11$.

En el diseño de un código hay dos objetivos prácticos opuestos:

1. Incorporar al código alguna información redundante de manera que tenga cierta robustez frente a errores debidos a su transmisión por un canal poco fiable.
2. Conseguir que el código no tenga ninguna información redundante para que el almacenamiento de la cadena de bits requiera el menor espacio posible.

Las asignaturas de teoría de códigos de los grados de matemáticas se suelen centrar en la primera situación, mientras que aquí nos centraremos en la segunda. En este contexto, hay un teorema bonito, práctico y sencillo (y posiblemente poco conocido en el ámbito matemático) que limita los milagros que podrían ocurrir al comprimir. Este resultado se suele conocer como *source coding theorem* y en [Sha48] (centrado en la comunicación) se considera en relación con la transmisión por un canal sin ruido, por ello también se denomina *noiseless coding theorem*. Con menos frecuencia se le llama *primer teorema de Shannon*.

Dado un código C y $s_i \in S$, la cadena de bits $C(s_i)$ tendrá cierta longitud $\ell(C(s_i))$. Por supuesto, para conseguir comprimir queremos que la longitud media de estas cadenas,

$$(4) \quad \ell(C) = \sum_i p_i \ell(C(s_i)),$$

sea pequeña. A partir de S podemos construir el conjunto \mathcal{S}_N de cadenas de N elementos. Las probabilidades p_i de los elementos de S inducen una medida de probabilidad en \mathcal{S}_N de manera natural suponiendo que los s_k que integran las cadenas se comportan como variables aleatorias independientes (esto no es cierto en muchas situaciones, por ejemplo en español detrás de una consonante es muy probable que vaya una vocal, volveremos sobre ello más adelante).

$$(5) \quad \mathcal{S}_N = \{s_{i_1} s_{i_2} \dots s_{i_N} : s_{i_k} \in S\}, \quad \text{Prob}(s_{i_1} s_{i_2} \dots s_{i_N}) = p_{i_1} p_{i_2} \dots p_{i_N}.$$

Para un código definido en \mathcal{S}_N , y no necesariamente definido en S , la longitud media es

$$(6) \quad \ell_N(C) = \sum_{i_1, \dots, i_N} p_{i_1} \dots p_{i_N} \ell(C(s_{i_1} \dots s_{i_N})).$$

El resultado que hemos anunciado es el siguiente:

Teorema (source coding theorem). Sea ℓ^* el mínimo de $\ell(C)$ sobre todos los posibles códigos definidos sobre S y H la entropía de S , entonces

$$(7) \quad H \leq \ell^* < H + 1.$$

Por supuesto, comprimir o codificar carácter a carácter podría ser absurdo si acumulan poca información, por ejemplo si son bits. Con un truco sencillo se pasa al caso general de cadenas.

Corolario. Sea ℓ_N^* el mínimo de $\ell_N(C)$ sobre todos los posibles códigos definidos en \mathcal{S}_N y H la entropía de S , entonces

$$(8) \quad H \leq \frac{\ell_N^*}{N} < H + \frac{1}{N}, \quad \text{en particular } \lim_{N \rightarrow \infty} \frac{\ell_N^*}{N} = H.$$

En breve, lo que asegura este resultado es que en media un fichero con N caracteres se puede comprimir como mucho hasta tamaño HN y que, de hecho, hay una igualdad asintótica. Es decir, que la entropía responde a la idea original de una medida de la cantidad de información por cada elemento de S .

Un ejemplo tonto para comprobar el resultado es que si $S = \{0, 1\}$ con probabilidades $p_1 = p_2 = 1/2$ estamos considerando ficheros al azar y no debería haber ninguna manera de comprimirlos, lo cual es coherente con que $H = 1$.

Pensemos ahora en ficheros binarios que típicamente tengan 75% de ceros y un 25% de unos. La entropía de S es $-\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} = 0.811$. En estas condiciones, un fichero de 100 Megas con un código adecuado se comprimiría en el límite hasta 81.1 Megas. Si codificamos los bits de uno en uno, obviamente no hay milagros y tendremos como poco los 100 Megas originales, pero si los agrupamos en bloques de 16 bits (dos bytes), el corolario nos asegura que existe un código tal que, en media, el tamaño en Megas estaría entre $100 \cdot 0.811$ y $100 \cdot (0.811 + 1/16) \approx 87.4$.

Pensemos ahora en ficheros compuestos por bytes, en los que el byte con valor 0 aparece la mitad de las veces y bytes con valores entre 1 y 20 con probabilidad 0.4 (esto es un modelo cualitativo de lo que sucede en el formato JPEG tras la cuantificación). La entropía es máxima si suponemos equidistribución en ausencia de más datos, es decir, que cada byte entre 1 y 20 tiene probabilidad $0.4/20$ y algo similar con los 235 números mayores de 20, entonces

$$(9) \quad -0.5 \log_2 0.5 - 20 \cdot \frac{0.4}{20} \log_2 \frac{0.4}{20} - 235 \cdot \frac{0.1}{235} \log_2 \frac{0.1}{235} = 3.877.$$

Esto significa que, en esta situación, cada byte lleva menos de 4 bits de información. Un fichero de 100 KB típicamente se podría comprimir con una codificación adecuada a $100 \cdot 1024 \cdot 3.877 \approx 397000$ bits que son aproximadamente 48.5 KB. Es decir, a menos de la mitad.

Para deducir el Corolario del Teorema, tomamos \mathcal{S}_N en vez de S , con \mathcal{S}_N todas las posibles cadenas de N elementos de S , de esta forma, ℓ^* se transforma en ℓ_N^* . La entropía correspondiente a \mathcal{S}_N es

$$(10) \quad - \sum_{i_1, i_2, \dots, i_N} p_{i_1} p_{i_2} \cdots p_{i_N} \log_2 (p_{i_1} p_{i_2} \cdots p_{i_N}) = - \sum_{i_1, i_2, \dots, i_N} p_{i_1} p_{i_2} \cdots p_{i_N} (\log_2 p_{i_1} + \cdots + \log_2 p_{i_N}).$$

Al desarrollar esta suma, el coeficiente de $-\log_2 p_1$ será

$$(11) \quad \sum_{k=1}^N \sum_{\substack{i_1, i_2, \dots, i_N \\ i_k=1}} p_{i_1} p_{i_2} \cdots p_{i_N} = N p_1 \sum_{j_1, \dots, j_{N-1}} p_{j_1} \cdots p_{j_{N-1}} = N p_1 \left(\sum_{i=1}^n p_i \right)^{N-1} = N p_1.$$

En general, el coeficiente de $-\log_2 p_i$ es $N p_i$ y la entropía de \mathcal{S}_N es NH con H la entropía de S . Así pues al aplicar (7) a \mathcal{S}_N se obtiene $NH \leq \ell_N^* < NH + 1$, que es (8).

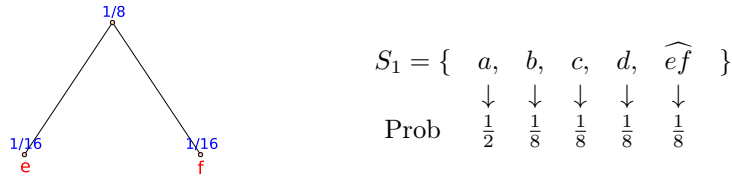
Codificación de Huffman

La prueba del *source coding theorem* es constructiva, es decir, en principio permitiría diseñar un código para el que ℓ_N^* satisfaga la desigualdad (8). En 1952, D.A. Huffman, siendo un estudiante de doctorado, creó un algoritmo que produce un código que minimiza $\ell(C)$. Lo más importante desde el punto de vista actual es que el algoritmo es sencillo y fácil de programar. En él se parte del conjunto $S_0 = S$ y se continúa hasta llegar a S_{n-1} , que sólo tiene un elemento, siendo el objetivo la creación de un árbol binario siguiendo las siguientes reglas:

1. Se toman los dos elementos de menor probabilidad de S_i y se añaden al árbol como hijos de un padre al que se le asigna la suma de las probabilidades.
2. Se construye S_{i+1} sustituyendo en S_i los hijos del paso anterior por su padre.

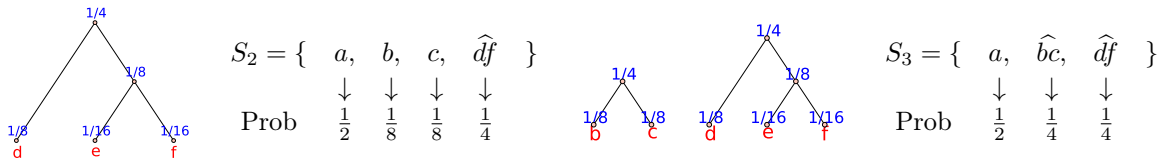
Dibujaremos el árbol de la manera habitual en teoría de grafos, con la raíz arriba y las hojas abajo, aunque esto es, por supuesto, irrelevante.

Como ejemplo, tomemos $S = \{a, b, c, d, e, f\}$ con probabilidades respectivas $p_1 = 1/2, p_2 = p_3 = p_4 = 1/8, p_5 = p_6 = 1/16$. Las probabilidades menores en $S_0 = S$ son las de e y f , por tanto el primer paso es

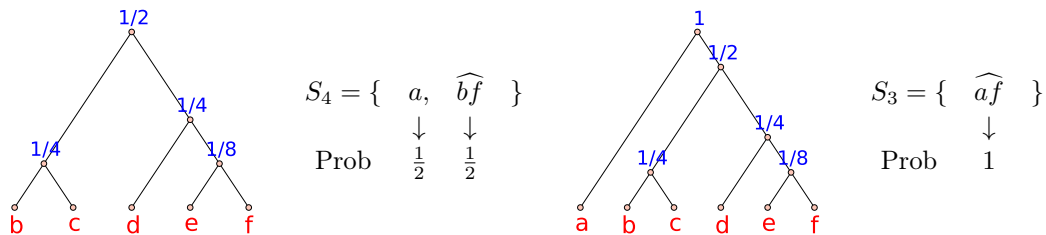


donde \widehat{ef} es el padre de e y f . En general, abreviaremos con \widehat{xy} el elemento del árbol que cumple que x es su primer hijo o desciende de él mientras que y es su segundo hijo o descendiente de él.

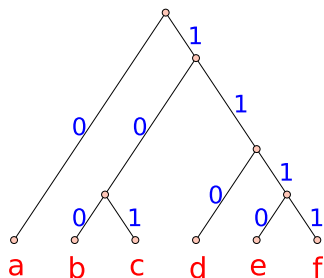
Ahora hay varias posibles elecciones de dos elementos de probabilidad $1/8$, la menor. Escojamos por ejemplo d y \widehat{ef} . En la siguiente iteración no hay ambigüedad y necesariamente hay que elegir b y c .



El siguiente paso será agrupar \widehat{bc} y \widehat{df} . Finalmente, el padre de ambos se agrupará con a y el padre de ellos completará la raíz del árbol.



Una vez construido el árbol, representaremos el camino desde la raíz a un elemento $s \in S$ en una de las hojas como una sucesión de ceros y unos, convencionalmente diremos que cero será el camino de la izquierda y uno el de la derecha. La cadena resultante es la que asignamos como codificación a s .



$a \rightarrow 000$ $d \rightarrow 110$
 $b \rightarrow 001$ $e \rightarrow 1110$
 $c \rightarrow 01$ $f \rightarrow 1111$

$$\ell(C) = 1 \cdot \frac{1}{2} + 3 \cdot \frac{1}{8} + 3 \cdot \frac{1}{8} + 3 \cdot \frac{1}{8} + 4 \cdot \frac{1}{16} + 4 \cdot \frac{1}{16} = \frac{17}{8}$$

$$H = -\frac{1}{2} \log_2 \frac{1}{2} - 3 \frac{1}{8} \log_2 \frac{1}{8} - 2 \frac{1}{16} \log_2 \frac{1}{16} = \frac{17}{8}$$

La longitud media es $\ell(C) = 17/8$ y coincide con la entropía; por tanto, según (7), es óptimo en cuanto a la longitud media.

La pregunta natural es cómo estamos seguros de que este algoritmo lleva a un código, es decir, que no hay cadenas distintas de elementos de S que den lugar a codificaciones iguales. La respuesta es más sencilla de lo que parece a primera vista y merece la pena ser generalizada.

Supongamos que tenemos una función $C : S \rightarrow \{\text{cadenas de bits}\}$ tal que para cualquier $b_1 b_2 \dots b_N \in \text{Im } C$, $b_i \in \{0, 1\}$ se tiene que la cadena $b_1 b_2 \dots b_k$ no está en $\text{Im } C$ para $1 \leq k < N$. Las funciones de este tipo se llaman *códigos prefijo* porque ningún $C(s_i)$ es el “prefijo” de otro (quizá el ejemplo más empleado en la práctica sea la codificación UTF-8). Es fácil ver que cualquiera de estas funciones es un código, porque una coincidencia de las imágenes de $s_{i_1} s_{i_2} \dots s_{i_k}$ y $s_{j_1} s_{j_2} \dots s_{j_l}$ implica que $C(s_{i_1})$ y $C(s_{j_1})$ violan la condición anterior si $s_{i_1} \neq s_{j_1}$, y si fueran iguales, basta considerar los siguientes.

Por la regla de asignar el 0 a la izquierda y el 1 a la derecha, la imagen de un código prefijo queda representado biyectivamente por un árbol binario (no completo en general, es decir, un padre puede tener un solo hijo) en el que las hojas corresponderán a los elementos que se codifican, porque si alguno estuviera en un nodo interior la condición de prefijo no se cumpliría. En particular, la codificación de Huffman es un código prefijo y por tanto un código válido.

La segunda cuestión es por qué la codificación de Huffman minimiza $\ell(C)$. Vamos a ver que lo hace entre el subconjunto de códigos prefijo. Como mencionaremos más adelante, en realidad siempre un código se puede transformar en código prefijo sin cambiar las longitudes, por tanto también minimiza entre todos los códigos.

Teorema. *El mínimo de $\ell(C)$ sobre todos los códigos prefijos se alcanza con la codificación de Huffman.*

Demostración. Digamos que $S = \{s_1, s_2, \dots, s_n\}$ con $p_1 \leq p_2 \leq \dots \leq p_n > 0$. Si $n > 1$, el árbol binario que corresponde al código prefijo óptimo lo podemos suponer completo porque si un padre no tuviera dos hijos, ese vértice se podría suprimir del árbol, reduciendo las longitudes de las codificaciones. La longitud de $C(s_1)$ debe ser máxima porque si $C(s_i)$ tuviera longitud mayor, $p_i \ell(C(s_i)) + p_1 \ell(C(s_1)) > p_i \ell(C(s_1)) + p_1 \ell(C(s_i))$ e intercambiar las codificaciones de s_1 y s_i sería ventajoso. De la misma forma, podemos cambiar las codificaciones de s_2 y la del hermano de s_1 en el árbol, porque $\ell(C(s_2))$ es máximo entre los $\ell(C(s_i))$, $i > 1$. Por tanto, podemos suponer que s_1 y s_2 corresponden a hojas hermanas en el árbol.

Ahora procedemos por inducción en n . Si $n = 2$ no hay nada que probar porque evidentemente la longitud media mínima es 1. Supongamos el resultado probado para $n - 1$ y sea $S' = \{z, s_3, s_4, \dots, s_n\}$ donde los s_i tienen las probabilidades de antes y z tiene probabilidad $p_1 + p_2$. Un código óptimo como antes en S induce un código en S' (quizá no óptimo) que en el árbol corresponde a sustituir las hojas de s_1 y s_2 por su padre z . Sean H y H' las codificaciones de Huffman en S y S' , entonces

$$(12) \quad \ell(C) = p_1 + p_2 + \ell(C') \quad \text{y} \quad \ell(H) = p_1 + p_2 + \ell(H').$$

La segunda igualdad es por la construcción del algoritmo de Huffman y la primera se sigue porque $p_1 \ell(C(s_1)) + p_2 \ell(C(s_2)) = (p_1 + p_2) \ell(C(s_1))$ se ha reemplazado por $(p_1 + p_2) \ell(C(z)) = (p_1 + p_2) (\ell(C(s_1)) - 1)$.

Por la hipótesis de inducción, $\ell(H') \leq \ell(C')$ y de (12) se deduce que $\ell(H) \leq \ell(C)$. □

La prueba del primer teorema de Shannon

Un punto fundamental en la deducción actual de (7) es la desigualdad de Kraft, introducida por L. Kraft en su tesis de máster de 1949 (un año después de [Sha48], a veces se ha dicho que Shannon no demostró sus teoremas) y probada independientemente por B. McMillan [McM56]. El siguiente enunciado incluye la propia desigualdad y su recíproco.

Teorema (Desigualdad de Kraft). *Dado un código prefijo C sobre $S = \{s_1, \dots, s_n\}$, las longitudes de las codificaciones $l_i = \ell(C(s_i))$ satisfacen*

$$(13) \quad \sum_{i=1}^n 2^{-l_i} \leq 1.$$

Además, dados enteros positivos l_i cualesquiera que verifiquen (13), siempre existe un código prefijo de forma que esos enteros son exactamente las longitudes de las codificaciones de los elementos de S .

La generalización de McMillan, aparte de añadir el recíproco, consiste en suprimir la hipótesis de que sea un código prefijo. De ello se deduce que dado un código C , existe siempre un código prefijo C_p tal que $\ell(C(s_i)) = \ell(C_p(s_i))$ y por ello, en cualquier resultado que hable de estas longitudes podemos suponer que el código es un código prefijo. De todas formas, los códigos prefijos son los más prácticos, porque coinciden con los *instantáneos*, los que pueden descodificarse al tiempo que se reciben [Rom92], sin esperar al final de la transmisión.

Demostración (del primer teorema de Shannon). Utilizando, por ejemplo, multiplicadores de Lagrange (o métodos más elementales), se tiene que la función $f(x_1, \dots, x_n) = -\sum p_i \log_2 x_i$ alcanza un mínimo en el *simplex* $\sum x_i = 1, 0 \leq x_i \leq 1$ cuando $x_i = p_i$.

Con la notación de la desigualdad de Kraft, tomemos $x_i = 2^{-l_i} / \sum 2^{-l_j}$, entonces la cota inferior de (7) se sigue de

$$(14) \quad H = f(p_1, p_2, \dots, p_n) \leq f(x_1, x_2, \dots, x_n) \leq f(2^{-l_1}, 2^{-l_2}, \dots, 2^{-l_n}) = \ell(C),$$

donde se ha aplicado (13) para la segunda desigualdad.

Para obtener la cota superior de (7), elijamos $l_i \in \mathbb{Z}^+$ tales que $-\log_2 p_i \leq l_i < 1 - \log_2 p_i$ (suponemos $p_i \neq 0$ sin pérdida de generalidad). Por el recíproco de la desigualdad de Kraft, debe existir un código prefijo con $l_i = \ell(C(s_i))$, ya que $\sum 2^{-l_i} \leq \sum p_i = 1$. Entonces $\ell(C) < \sum p_i(1 - \log_2 p_i) = 1 + H$ y se obtiene el resultado deseado. \square

Demostración (de la desigualdad de Kraft). Digamos que $l_1 \leq l_2 \leq \dots \leq l_r < l_{r+1} = \dots = l_n$. Consideremos el conjunto B de todas las cadenas de l_n bits. Si $i \leq r$, el conjunto $B \cap \text{Im } C$ no contiene a ninguna de las $2^{l_n - l_i}$ cadenas que comienzan como $C(s_i)$, mientras que $C(s_i) \in B$ para $i > r$. Por tanto

$$(15) \quad |B| = |B \cap \text{Im } C| + |B - B \cap \text{Im } C| \geq n - r + \sum_{i=1}^r 2^{l_n - l_i} = \sum_{i=1}^n 2^{l_n - l_i}$$

y basta notar que $|B| = 2^{l_n}$. Otra forma de proceder más simple es asociar a B el árbol binario completo. A las $2^{l_n - l_j}$ cadenas de B que empiezan por $C(s_j)$ les corresponden las hojas de un subárbol binario completo. La condición de prefijo asegura que estos subárboles son disjuntos y se deduce la desigualdad.

Para el recíproco consideramos el árbol binario completo T de altura (*profundidad*) l_n , el cual tiene 2^{l_n} hojas. Establecemos el siguiente algoritmo: Para cada $i \leq r$ consideramos un vértice que diste l_i de la raíz (que tenga profundidad l_i) y eliminamos todos sus descendientes, entre ellos $2^{l_n - l_i}$ hojas de T . Si queremos que este algoritmo pueda llevarse a cabo y que al final queden al menos $n - r$ hojas de T , necesitamos $2^{l_n} - \sum_{i=1}^r 2^{l_n - l_i} \geq n - r$, lo cual equivale a la hipótesis (13). Tras el algoritmo, si quedasen más de $n - r$ hojas de T , eliminamos las sobrantes. El resultado será un árbol T' cuyas hojas tienen profundidades exactamente los l_i . Así pues, T' corresponde a un código prefijo con las propiedades esperadas. \square

Métodos de diccionario

Imaginemos un tipo de ficheros en que cada byte tiene un 99% de posibilidades de ser igual al anterior y en el 1% restante cambia aleatoriamente. Si el fichero es suficientemente grande, la probabilidad de aparición de cada posible byte será 1/256 y por tanto $H = 8$, que indica que cada byte tiene 8 bits. Si aplicamos (8), se deduce que no hay compresión posible. Sin embargo está claro que hay un método de compresión muy sencillo, llamado RLE (*run-length encoding*): Por cada cadena de bytes iguales usamos un byte para indicar el número de repeticiones y otro para indicar el byte que se repite. En caso de que haya más de 255 repeticiones (lo cual es poco probable) se puede dividir la cadena en dos partes. Por ejemplo, 200A100A100B significa 300 letras A seguidas de 100 letras B. Como típicamente cada cadena de bytes iguales se repite 100 veces, la compresión es del orden de 1/50.

La paradoja proviene de que en (8) las probabilidades se asignan a través de (5) suponiendo la independencia de los caracteres, lo cual no es el caso aquí.

La codificación de Huffman tomando como S todas las posibles cadenas de 100 bytes que aparecen en el fichero, en principio es muy ventajosa aquí pero ese dato viene de un conocimiento previo del tipo de fichero. Por otro lado, S puede llegar a ser bastante grande y el árbol correspondiente complicado. Esto muestra dos de las dificultades para aplicar la codificación de Huffman como compresor de propósito general:

1. Depende de la división en símbolos de S que requiere una estadística previa costosa.
2. El árbol debe ser conocido por el descompresor y almacenarlo puede ser poco eficiente.

A partir de los años 70 del siglo XX, se diseñaron algunos métodos, a veces conocidos genéricamente como Lempel-Ziv, por los autores del trabajo original [ZL77], que tenían en mente especialmente textos (los ficheros multimedia no estaban apenas desarrollados entonces) pero con la aspiración de ser algoritmos universales, como dice el título de [ZL77] y justifican sus resultados teóricos. Los tres más conocidos son LZ77, LZ78 y LZW. Son parecidos pero quizá a través de LZ78 sea más sencillo entender la idea. Lo veremos con un poco más de detalle e indicaremos en que consisten las variantes para los otros.

El algoritmo LZ78. Para simplificar, suponemos que al final de la cadena que queremos comprimir aparece un carácter de *fin de fichero* (EOF), que representamos con # y que no está en otro lugar.

Lo primero es descomponer la cadena en *frases*, donde una frase es una subcadena que añade un carácter nuevo a otra que haya aparecido antes. Por ejemplo, las divisiones en frases de `salsa_salada` (donde `␣` indica un espacio) y de `rellena_la_encuesta`, serían

`|s|a|l|sa|␣|sa|l|a|d|a|#` y `|r|e|l|l|e|n|a|␣|l|a|␣e|n|c|u|e|s|t|a|#`

Por supuesto cada frase aparece sólo una vez. La división en frases se corresponde a un diccionario (una lista), cuya entrada k es la k -ésima frase y por convenio diremos que la entrada 0 es la cadena vacía. En los ejemplos anteriores los diccionarios serían

0	∅	5	␣	y	0	∅	5	n	10	nc
1	s	6	sal		1	r	6	a	11	u
2	a	7	ad		2	e	7	␣	12	es
3	l	8	a#		3	l	8	la	13	t
4	sa				4	le	9	␣e	14	a#

El propósito del carácter de fin de fichero es asegurarnos de que la última frase no se quede sin almacenar en el diccionario.

El texto original se codifica asignando a cada frase un par, dado por su lugar en el diccionario al quitarle el último carácter y por ese último carácter. En los ejemplos anteriores se tendrían las codificaciones

(0, s), (0, a), (0, l), (1, a), (0, ␣), (4, l), (2, d), (2, #)

y

(0, r), (0, e), (0, l), (3, e), (0, n), (0, a), (0, ␣), (3, a), (7, e), (5, c), (0, u), (2, s), (0, t), (6, #)

Es fácil ver que a partir de estas codificaciones se pueden recuperar los diccionarios y con ellos, descodificar. Por ejemplo, en el primer caso, el descodificador lee $(0, \mathbf{s})$ y sabe que tiene que poner una \mathbf{s} en la posición 1, de la misma forma, las posiciones 2 y 3 pasarán a estar ocupados por \mathbf{a} y \mathbf{l} . A continuación se lee $(1, \mathbf{a})$, y esto significa que hay que poner en la siguiente posición, la 4, lo que hay en la posición 1 y añadirle una \mathbf{a} . En general, cuando se lee un par (k, x) , se añade al final del diccionario lo que hay en la posición k terminado en x . El texto original es simplemente el diccionario leído de principio a fin.

El primer texto en ASCII ocuparía 13 bytes (13 caracteres) y si en la codificación reservamos un byte para cada número (esto no es posible para diccionarios de más de 256 elementos), entonces la “compresión” ocupa 16 bytes. Esto es más que el original. El segundo ejemplo es todavía peor, originalmente tiene 20 bytes y su codificación tiene 28 bytes. Esta aparente ineficiencia del algoritmo es una falsa impresión que se debe a que el texto es muy breve. Cuando es suficientemente grande, el número de coincidencias con elementos del diccionario aumenta rápidamente. Por ejemplo, los 1869610 caracteres del texto que habíamos manejado antes, dan lugar a un diccionario de tamaño 281192. Es decir, la codificación requiere esta cantidad de pares. El número que aparece en cada par se puede representar con 19 bits, porque $2^{19} > 281192$ y el carácter con 8, en total 27 bits por cada par, lo que hacen 949023 bytes, lo cual rebaja la cota de entropía de los caracteres pero está por encima de lo que consiguen los compresores habituales.

En la práctica, tener que manejar diccionarios grandes en los procesos de codificación y descodificación puede ser ineficiente sin métodos adecuados para gestionarlos [Sal02].

Universalidad del algoritmo LZ78. La fuerza de la compresión con LZ78 está ligada al número de frases F que se guardan en el diccionario. Cada par de la codificación emplea un carácter y una referencia a un frase del diccionario que requerirá $\lceil \log_2 F \rceil$ bits (donde $\lceil x \rceil$ es el menor entero $n \geq x$). Como hay F pares, la longitud en bits tras aplicar el algoritmo a una cadena s será

$$(16) \quad \ell_{LZ}(s) = F(\lceil \log_2 F \rceil + b)$$

donde b es el número de bits necesarios para codificar los caracteres que componen s .

Para comparar LZ78 con el mejor código involucrado en (8), consideramos la longitud media al codificar elementos de \mathcal{S}_N , las cadenas de longitud N , con la probabilidad indicada en (5), es decir,

$$(17) \quad \ell_{LZ}^*(N) = \sum_{i_1, \dots, i_N} p_{i_1} p_{i_2} \cdots p_{i_N} \ell_{LZ}(s_{i_1} s_{i_2} \cdots s_{i_N}).$$

Lo que deseamos es que esta longitud media sea pequeña.

Teorema. *Con la notación anterior se tiene*

$$(18) \quad \lim_{N \rightarrow \infty} \frac{\ell_{LZ}^*(N)}{N} = H$$

donde H es la entropía de S .

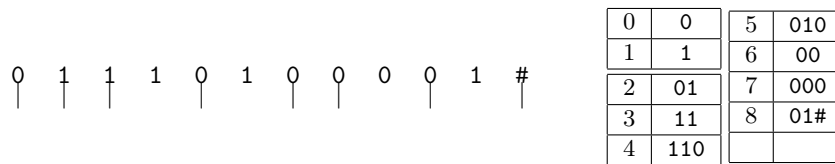
Esto prueba que LZ78 es *universal*, es decir, que asintóticamente comprime tanto como puede comprimir, sea cuales sean las probabilidades en S con una entropía fijada. De nuevo, es importante recordar que la definición de $\ell_{LZ}^*(N)$ implícitamente supone el modelo de variables aleatorias independientes de (5). este modelo puede ser muy poco realista en aplicaciones concretas (por ejemplo en ficheros de texto).

El algoritmo LZW. Es una variante de LZ78 que carga un diccionario inicial con todos los posibles caracteres individuales (bytes, elementos de S). Tiene como ventaja que permite eliminar el segundo elemento de cada par de LZ78 a cambio de una codificación un poco más compleja y e alguna excepción en la descodificación. La división en frases es similar a la de LZ78 pero con dos particularidades:

1. Los caracteres individuales están precargados en el diccionario y por tanto no dan lugar a frases.
2. Cada carácter del final de una frase se incluye en el principio de la siguiente.

Todas las frases se almacenan en el diccionario ordenadamente (como en LZ78), tras las entradas iniciales, y la codificación es el lugar que ocupa en el diccionario cada frase sin su último carácter. Suponemos, como antes, que al final de los datos hay un carácter especial para que la última frase se incorpore al diccionario. En este caso, este carácter no aparecerá en la descodificación.

Para simplificar, veamos un ejemplo bit a bit (lo cual no es realista). De esta forma el diccionario inicial contiene simplemente 0 y 1. Tomemos 01110100001# como la cadena a codificar. La división en frases y el diccionario correspondiente serían:



Se han colocado las barras divisorias de frases debajo para resaltar que cada carácter sobre una frase está tanto en la frase de la izquierda como en la de la derecha. Para codificar, basta poner la referencia de todas las frases sin su último carácter. Esto es lo mismo que si en la división en frases (sin #) los caracteres fronterizos pasasen a la derecha y buscásemos los resultados en el diccionario

$$| 0 | 1 | 11 | 01 | 0 | 00 | 01 | \longrightarrow 0, 1, 3, 2, 0, 6, 2.$$

La descodificación merece unos comentarios adicionales. En principio se aplica el mismo esquema que en LZ78: con el número en curso de la codificación y el anterior, se crea un nuevo elemento del diccionario a partir de los anteriores, y una vez recuperado el diccionario completo basta leer los elementos de la parte no precargada, omitiendo los últimos caracteres de cada uno de ellos para que no se repitan. El problema es que a veces uno de los códigos se refiere a un elemento del diccionario que no ha sido creado. Eso sólo puede ocurrir si la descodificación buscada coincide con la anterior añadiendo su primer carácter al final.

En el ejemplo anterior, el descodificador partiría del diccionario inicial $0 \rightarrow 0$, $1 \rightarrow 1$ y al leer 0 y 1 en la codificación generaría 01 y ampliaría el diccionario con $2 \rightarrow 01$. Después leería 3, esta parte del diccionario no está todavía creada. Como la última descodificación fue 1, el 3 descodifica como 11 (primer carácter al final) y se amplía el diccionario con $3 \rightarrow 11$. La descodificación entonces sería por ahora 0111. Después se leería 2 que corresponde a 01 y se añadiría al diccionario la nueva frase $4 \rightarrow 110$. Siempre cada frase añadida al diccionario es la última descodificación más el primer carácter de la que está en curso. El esquema completo para el ejemplo es:

Código	0	1	3*	2	0	6*	2
Descodificación	0	1	11	01	0	00	01
Diccionario		2→01	3→11	4→110	5→010	6→00	7→000

Los casos que llevan * son los que se refieren a elementos del diccionario que no han sido creados

Quizá es más sencillo entender el proceso con letras y cuando no hay que aplicar la regla “especial” (realmente en los ficheros que se comprimen bien, la situación especial es común). Digamos que partimos de un diccionario inicial $0 \rightarrow \square$, $1 \rightarrow h$, $2 \rightarrow i$, $3 \rightarrow o$, $4 \rightarrow p$, y queremos descodificar 1, 2, 4, 0, 5, 7, 9, 3. El esquema sería:

Código	1	2	4	0	5	7	9	3
Descodificación	h	i	p	□	hi	p□	hip	o
Diccionario		5→hi	6→ip	7→p□	8→□h	9→hip	10→p□h	11→hipo

Es decir, $hip_{\square}hip_{\square}hipo_{\#}$.

El algoritmo LZ77. Éste es el algoritmo original que dio lugar al resto. Más que en diccionarios se basa en el concepto de *ventana móvil* (*sliding window*) que esencialmente restringe la búsqueda de coincidencias a una porción de los datos de longitud predeterminada. La ventaja es que el número de bits para indicar referencias se conoce de antemano y no puede crecer. Por otro lado, la codificación emplea ternas, en lugar de los pares de LZ78 o los números sueltos de LZW.

No entraremos en detalles (que pueden consultarse en el original [ZL77], aunque es un poco teórico, o en [Sal02]), simplemente ilustraremos la situación con un ejemplo. Codifiquemos $sals_{\square}salada$ con LZ77 usando una ventana móvil de tamaño 10. Esta ventana tiene dos partes, una corresponde al “pasado” y otra al “futuro”. Digamos que sus tamaños son respectivamente 6 y 4 (según [Sal02], en la práctica la primera es de miles de bytes y la segunda de decenas de bytes).

En cada paso se consideran las cadenas que parten del primer carácter futuro y se busca una coincidencia lo más larga posible con cadenas que empiezan en el pasado. La codificación se expresa con una terna (a, b, c) donde a dice la distancia a la que está la coincidencia, b su longitud y c es el primer carácter futuro en el que no hay coincidencia, hacia el que se moverá la ventana de futuro en el siguiente paso.

Inicialmente no hay nada con que comparar y entonces las primeras codificaciones son triviales con $a = b = 0$:

□	□	□	□	□	□	s	a	l	s	→ (0, 0, s)
□	□	□	□	□	s	a	l	s	a	→ (0, 0, a)
□	□	□	□	s	a	l	s	a	□	→ (0, 0, 1)

Una vez que la segunda s pasa a ser el carácter futuro inicial, se produce una coincidencia de longitud 2, la cadena sa , y hay que ir 3 cuadros hacia el pasado para encontrarla:

□	□	□	s	a	l	s	a	□	s	→ (3, 2, □)
---	---	---	---	---	---	---	---	---	---	-------------

El resto de la codificación, hasta llegar al fin de los datos sería:

□	s	a	l	s	a	□	s	a	l	→ (0, 0, s)
s	a	l	s	a	□	s	a	l	a	→ (6, 3, a)
a	□	s	a	l	a	d	a	□	□	→ (0, 0, a)
□	s	a	l	a	d	a	□	□	□	→ (4, 1, #)

De nuevo, en este ejemplo, la codificación aumenta mucho lo que ocuparían habitualmente los datos iniciales. Esta impresión viene de que es un texto muy breve.

Una situación no ilustrada en el ejemplo anterior, es que no se prohíbe que una coincidencia se extienda hasta los caracteres futuros. Por ejemplo $0\dots 01\|110\dots$ se codificaría con $(1, 3, 0)$.

Referencias

- [Mac03] D. J. C. MacKay. *Information theory, inference and learning algorithms*. Cambridge University Press, New York, 2003.
- [McM56] B. McMillan. Two inequalities implied by unique decipherability. *IEEE Trans. Information Theory*, 2(4):115–116, 1956.

- [Rén84] A. Rényi. *A diary on information theory*. Akadémiai Kiadó (Publishing House of the Hungarian Academy of Sciences), Budapest, 1984. With a foreword by Pál Révész, Translated from the Hungarian by Zsuzsanna Makkai-Bencsáth.
- [Rom92] S. Roman. *Coding and information theory*, volume 134 of *Graduate Texts in Mathematics*. Springer-Verlag, New York, 1992.
- [Sal02] D. Salomon. *A guide to data compression methods*. Springer-Verlag, 2002.
- [Sha48] C. E. Shannon. A mathematical theory of communication. *Bell System Tech. J.*, 27:379–423, 623–656, 1948.
- [ZL77] J. Ziv and A. Lempel. A universal algorithm for sequential data compression. *IEEE Trans. Information Theory*, IT-23(3):337–343, 1977.