

3.5. CODIFICACIÓN Y ENTROPÍA

Codificar un conjunto de caracteres que aparecen en un "texto" es asignar a cada uno de ellos un código, generalmente binario (0 y 1). El objetivo es minimizar el número de bits necesarios para codificar el texto.

"Texto" = X ; los símbolos de X se escriben $S = \{x_1, \dots, x_N\}$

Ejemplo 1. X = página de un libro, S = letras del alfabeto, incluidos los signos de puntuación (, , . , ;) y los espacios.

Ejemplo 2. X = fotografía digitalizada en (R, G, B)
 $S = \{0, 1, \dots, 255\}$ en cada canal

Ejercicio 3.5.1. $S = \{A, O, E, S\}$, $X = AASAEEOA$

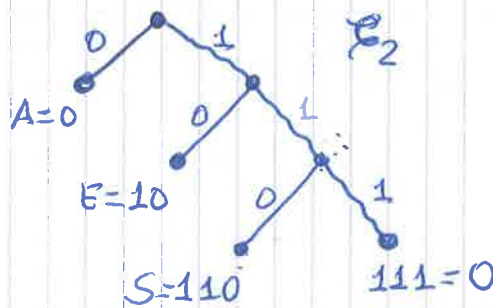
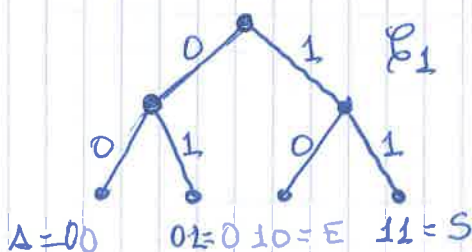
Ejemplo de Código

$$\mathcal{C}_1 = \{A=00, O=01, E=10, S=11\}$$

$$AASAEEOA \rightarrow 0000110010100001 \quad (16 \text{ bits})$$

$$\mathcal{C}_2 = \{A=0, E=10, S=110, O=111\}$$

$$AASAEEOA \rightarrow 00110010100111 \quad (14 \text{ bits})$$



l_k = longitud del código binario que representa a x_k

$P_k = \Pr(x_k \in X)$ la probabilidad de que x_k esté en X

La media de bits necesarios para codificar cada símbolo en un código \mathcal{C} es $M_X(\mathcal{C}) = \sum_{k=1}^N l_k P_k$.

El objetivo es encontrar un código \mathcal{C} para el texto X que haga mínima $M_X(\mathcal{C})$

Ejercicio 3.5.2 $X = AASAEEAO, S = \{A, O, E, S\}$

$Pr(A \in X) = \frac{4}{8} = \frac{1}{2}; Pr(O \in X) = \frac{1}{8}, Pr(E \in X) = \frac{1}{4}, Pr(S \in X) = \frac{1}{8}$

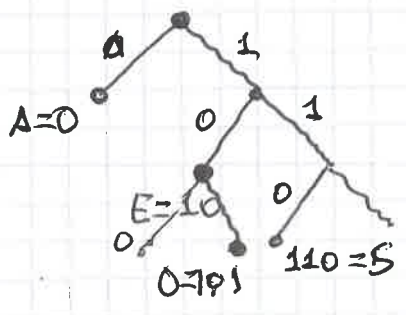
$M_X(\mathcal{C}_1) = 2 \times \frac{1}{2} + 2 \times \frac{1}{8} + 2 \times \frac{1}{4} + 2 \times \frac{1}{8} = 2$

$M_X(\mathcal{C}_2) = 1 \times \frac{1}{2} + 3 \times \frac{1}{8} + 2 \times \frac{1}{4} + 3 \times \frac{1}{8} = 1 + \frac{6}{8} = 1 + \frac{3}{4} = 1.75$

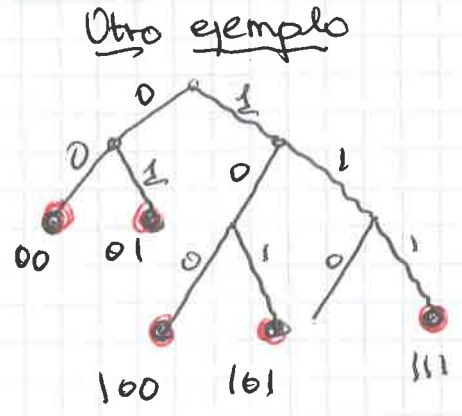
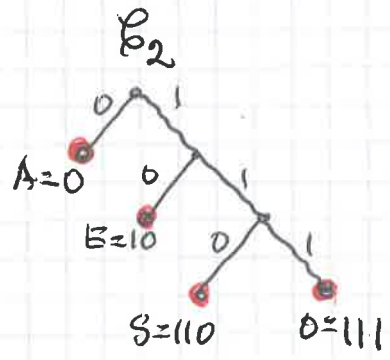
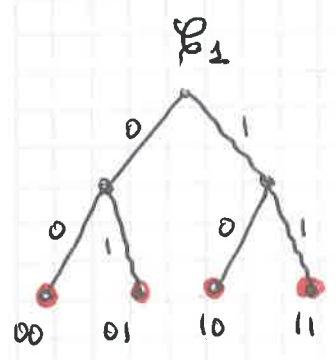
Los códigos con palabras binarias de longitud variable, NO siempre son buenos para decodificar. Por ejemplo,

$\mathcal{C}_3 = \{A=0, E=10, S=110, O=101\}$

no es bueno porque 1010 pueda ser EE ó OA



Def 3.5.1. Un código \mathcal{C} cumple la condición del prefijo si no hay códigos de \mathcal{C} que sean el comienzo de otro código de \mathcal{C} .



Ejemplo de códigos que cumplen la condición del prefijo

3.5.1. LA ENTROPÍA DE SHANNON

Def 3.5.2. Dado un texto X en el que cada símbolo X_k tiene probabilidad $p_k = \Pr(X_k \in X)$, la entropía de X es

$$\mathcal{E}(X) = \sum_{k=1}^N p_k \log_2 \frac{1}{p_k} \geq 0$$

- Supongamos que $p_1 = \Pr(X_1 \in X) = 1$ y $p_k = \Pr(X_k \in X) = 0$ si $k=2, \dots, N$. Entonces $\mathcal{E}(X) = 1 \cdot \log_2 \frac{1}{1} + \sum_{k=2}^N 0 \cdot \log_2 \frac{1}{0} = 1$ pg. $\lim_{x \rightarrow 0^+} x \log_2 \frac{1}{x} \stackrel{\text{L'Hôpital}}{=} 0$

- Supongamos que $p_k = \Pr(X_k \in X) = \frac{1}{N}$, $k=1, \dots, N$. Entonces $\mathcal{E}(X) = \sum_{k=1}^N \frac{1}{N} \log_2 \frac{1}{\frac{1}{N}} = (\log_2 N) \frac{1}{N} \sum_{k=1}^N 1 = \log_2 N$

Ejercicio 3.5.3 $X = \{AASAEEAO\}$, $\mathcal{S} = \{A, O, E, S\}$

Calcula $\mathcal{E}(X)$

$$S/ \Pr(A \in X) = \frac{1}{2}, \Pr(O \in X) = \frac{1}{8}, \Pr(E \in X) = \frac{1}{4}, \Pr(S \in X) = \frac{1}{8}$$

$$\begin{aligned} \mathcal{E}(X) &= \frac{1}{2} \log_2 2 + \frac{1}{8} \log_2 8 + \frac{1}{4} \log_2 4 + \frac{1}{8} \log_2 8 \\ &= \frac{1}{2} \cdot 1 + \frac{1}{8} \cdot 3 + \frac{1}{4} \cdot 2 + \frac{1}{8} \cdot 3 = 1.75 \end{aligned}$$

Proposición 3.5.3 Para cualquier texto X escrito con N símbolos, $0 \leq \mathcal{E}(X) \leq \log_2 N$.

D/ Halla el máximo de $\mathcal{E}(X) = \sum_{k=1}^N p_k \log_2 \frac{1}{p_k} = - \sum_{k=1}^N p_k \log_2 p_k$ con las condiciones $0 \leq p_k \leq 1$ y

$$\sum_{k=1}^N p_k = 1. \text{ Lagrange}$$

$$F(p_1, \dots, p_N, \lambda) = - \sum_{k=1}^N p_k \log_2 p_k + \lambda \left(\sum_{k=1}^N p_k - 1 \right)$$

$$\left. \begin{aligned} \frac{\partial F}{\partial P_k} &= -\log_2 P_k - \frac{1}{\ln 2} + \lambda = 0 \\ \frac{\partial F}{\partial \lambda} &= \sum_{k=1}^N P_k - 1 = 0 \end{aligned} \right\} \Rightarrow \log_2 P_k = \lambda - \frac{1}{\ln 2}$$

$$P_k = 2^\lambda \cdot 2^{-\frac{1}{\ln 2}} = 2^\lambda \cdot 2^{-\log_2 e}$$

$$\Rightarrow P_k = \frac{2^\lambda}{e}$$

$$1 = \sum_{k=1}^N P_k = \sum_{k=1}^N \frac{2^\lambda}{e} = \frac{2^\lambda}{e} N \Rightarrow 2^\lambda = \frac{e}{N} \Rightarrow \lambda = \log_2 \frac{e}{N}$$

$$P_k = \frac{2^\lambda}{e} = \frac{2^{\log_2 \frac{e}{N}}}{e} = \frac{\frac{e}{N}}{e} = \frac{1}{N} \quad \text{Entonces}$$

$$\mathcal{L}(X) = \sum_{k=1}^N \frac{1}{N} \log_2 N = \log_2 N.$$

Teorema 3.5.4 (C. Shannon)

$X =$ texto con símbolos $S = \{x_1, \dots, x_N\}$; $P_k = \Pr(x_k \in X)$, $k=1, \dots, N$; \mathcal{L} código binario con la condición del prefijo.

Entonces,

$$M_X(\mathcal{L}) \geq \mathcal{L}(X) = \sum_{k=1}^N P_k \log_2 \frac{1}{P_k}$$

Además, existe un código binario \mathcal{L}' con la condición del prefijo tal que

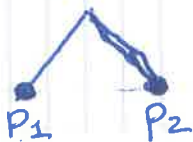
$$M_X(\mathcal{L}') \leq \mathcal{L}(X) + 1.$$

3.5.2. EL CÓDIGO DE HUFFMAN

Este algoritmo permite hacer un código binario con la condición del prefijo que minimiza $M_X(\mathcal{L})$

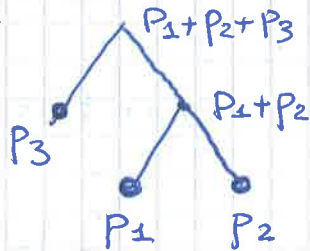
1. $P(x) = \Pr(x \in X)$. Ordenar de menor a mayor las probabilidades $P(x)$. Supongamos $P_1 \leq P_2 \leq \dots \leq P_N$

2. Tomar los símbolos con los dos probabilidades menores y asignarles a dos hijos de un árbol binario

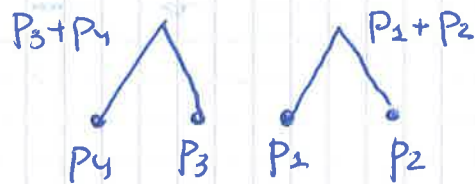


3. Calcular $P_1 + P_2$ y ordenar la lista $P_1 + P_2, P_3, P_4, \dots, P_N$ de menor a mayor

4'. Si tenemos $P_1 + P_2 \leq P_3 \leq P_4 \leq \dots \leq P_N$ o' $P_3 \leq P_1 + P_2 \leq P_4 \leq \dots \leq P_N$ hacemos el árbol



4''. Si tenemos $P_3 \leq P_4 \leq \dots \leq P_N$ hacemos el árbol



5. Ordenar la nueva lista de probabilidades según cada caso y continuar hasta asignar la última probabilidad a una hoja del árbol.

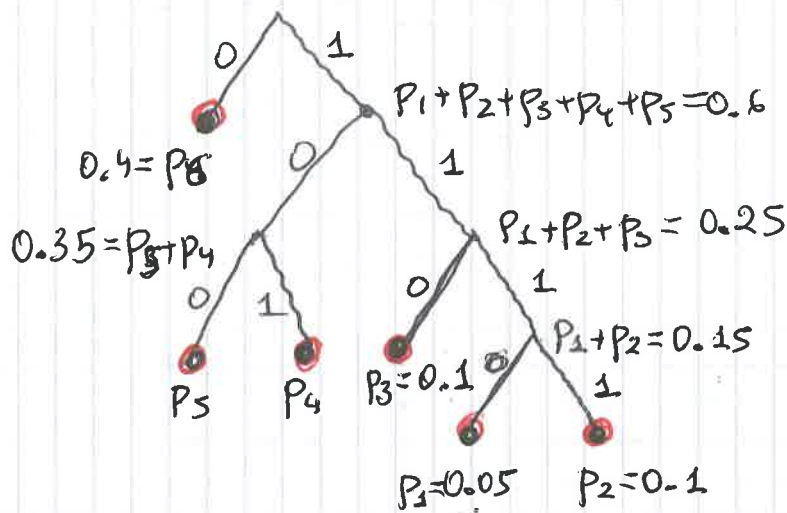
Ejercicio 3.5.4 a) Construir el código de Huffman \mathcal{C} para los símbolos $\{a, b, c, d, e, f\}$ de los que conocemos que aparecen en un texto X con probabilidades

$$p(a) = 0.2, \quad p(b) = 0.05, \quad p(c) = 0.15, \quad p(d) = 0.4$$

$$p(e) = 0.1, \quad p(f) = 0.1$$

b) Calcular $\mathcal{E}(X)$ y $M_X(\mathcal{C})$

(a) $P_1 \leq P_2 \leq P_3 \leq P_4 \leq P_5 \leq P_6$
 0.05 0.1 0.1 0.15 0.2 0.4
 p(b) p(e) p(f) p(c) p(a) p(d)



$$\left\{ \begin{array}{l} \boxed{P_3 \leq P_1+P_2} \leq P_4 \leq P_5 \leq P_6 \\ 0.1 \quad 0.15 \quad 0.15 \quad 0.2 \quad 0.4 \end{array} \right\} \text{ Caso 4'}$$

$$\left\{ \begin{array}{l} P_4 \leq P_5 \leq P_1+P_2+P_3 \leq P_6 \\ 0.15 \quad 0.2 \quad 0.25 \quad 0.4 \end{array} \right\} \text{ Caso 4''}$$

$$\underbrace{P_1+P_2+P_3=0.25} \leq P_5+P_4=0.35, \leq P_6=0.4 \quad \text{Caso 4'}$$

$$\left. \begin{array}{l} b \rightarrow 1110, \quad e \rightarrow 1111, \quad f \rightarrow 110, \quad c \rightarrow 101, \\ a \rightarrow 100, \quad d \rightarrow 0 \end{array} \right\}$$

$$b) \quad \mathcal{E}(X) = \sum_{k=1}^6 P_k \log_2 \frac{1}{P_k} = 0.05 \log_2 \frac{1}{0.05} + (0.1 \log_2 \frac{1}{0.1}) \cdot 2 + 0.15 \log_2 \frac{1}{0.15} + 0.2 \log_2 \frac{1}{0.2} + 0.4 \log_2 \frac{1}{0.4} = 2.2841$$

$$M_X(\mathcal{E}) = (0.2) \cdot 3 + 4 \cdot (0.05) + 3 \cdot (0.15) + 1 \cdot (0.4) + 4 \cdot (0.1) + 3 \cdot (0.1) = 2.35 \approx \mathcal{E}(X)$$

Ejercicio 3.5.5 (Para entregar el ^{24/07} 23/07, Sexta feira)

(a) Construir el código de Huffman \mathcal{C} para el texto

$$X = 13524 135 13524 13513524 111$$

hecho con los símbolos de $S = \{1, 2, 3, 4, 5\}$

b) Calcular $\mathcal{E}(X)$ y $M_X(\mathcal{C})$, comprobando que $M_X(\mathcal{C}) \geq \mathcal{E}(X)$
