

Huffman Code

jueves, 9 de noviembre de 2023 09:04

- Σ alfabeto de caracteres $\{a, b, c, d, e, f\}$
- F archivo que contiene 100k caracteres.
- Conocemos la frecuencia que aparecen los caracteres en F

	a	b	c	d	e	f
f (miles)	45	13	12	16	9	5

- Queremos almacenar F
- Objetivo, usar el menor espacio posible.

Tools:

Código binario de caracteres: (o código)

$$\mathcal{C}: \Sigma \longrightarrow \{0,1\}^+ \quad \text{inyectiva}$$

"x" \longmapsto $\mathcal{C}(x)$

Decimos que $\mathcal{C}(x)$ es el código del carácter 'x'.

	a	b	c	d	e	f
Frequency (in thousands)	45	13	12	16	9	5
Fixed-length codeword	000	001	010	011	100	101
Variable-length codeword	0	101	100	111	1101	1100

~~10~~ ~~11~~ ~~010~~ ~~000~~ ~~0111~~ ~~00111~~

000001 ✓

cab
acb
aabb



$\mathcal{C}(x)$
 $\mathcal{C}(y)$

Código libre de prefijos

(Códigos donde $x \neq y \Rightarrow \mathcal{C}(x)$ no es prefijo de $\mathcal{C}(y)$)

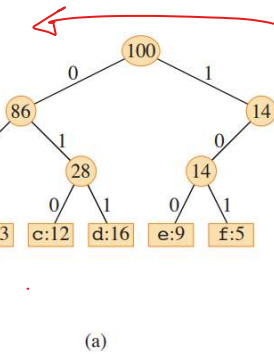
Problema: Obtener m

$$\min \left\{ \sum_{x \in F} |\mathcal{C}(x)| : \mathcal{C} \text{ código libre de prefijos de } \Sigma \right\}$$

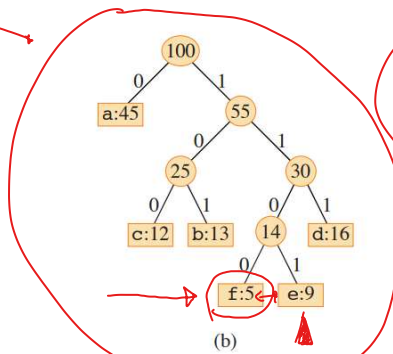
Prefix-free codeBinary Tree, leaves with characters

abbcddedf
 1 1 3 3 3 3 4 4

	f	d
a	1	1
b	2	3
c	1	3
d	2	3
e	1	4
f	1	4



(a)



(b)

$|C(b)|$

3

1

Prop: optno se alcanza con un full binary tree.

Propiedad: $d_T: \Sigma \rightarrow \mathbb{Z}_0^+$
 $x \mapsto d_T(x)$ profundidad

$$\text{val}(T) = \sum_{x \in \Sigma} d_T(x) \cdot f(x)$$

Prob: obtener mn:

$\{ \text{val}(T) : T \text{ tree associated to a pfc } \subseteq \Sigma \}$

Lemma: Sea $c \in \Sigma$ el caracter con menor frecuencia

$\Rightarrow \exists$ un T optno T_0 en el que c se ubica en un nodo con profundidad máxima

Sea T un óptimo:

Caso 1: c tiene prof máxima - Trivial

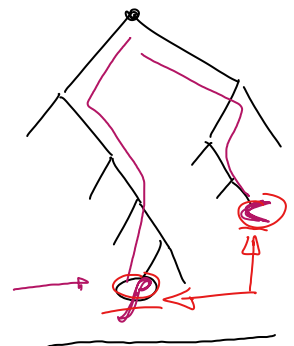
Caso 2: c no tiene prof máxima -

$\Rightarrow \exists p \in \Sigma$ con prof máxima

$$d_T(p) > d_T(c)$$

Construimos T' swapando c con p .

Vemos que



$$d_{T'}(x) = \begin{cases} d_T(x) & x \notin \{c, p\} \\ d_T(p) & x = c \\ d_T(c) & x = p \end{cases}$$

Comparando T y T' :

$$\begin{aligned} \underline{\text{val}(T) - \text{val}(T')} &= \sum_{x \in \Sigma} f(x) \cdot d_T(x) - \sum_{x \in \Sigma} f(x) \cdot d_{T'}(x) \\ &= f(c) \cdot d_T(c) + f(p) \cdot d_T(p) - f(c) \cdot d_{T'}(c) - f(p) \cdot d_{T'}(p) \\ &= f(c) \cdot d_T(c) + f(p) \cdot d_T(p) - f(c) \cdot d_T(p) - f(p) \cdot d_T(c) \\ &= [f(c) - f(p)] [d_T(c) - d_T(p)] \geq 0. \end{aligned}$$

$\text{val}(T) \geq \text{val}(T')$. $\rightarrow T'$ es un árbol óptimo

Lema 2: • Sean b y c los caracteres con menor frecuencia ($f(c) \leq f(b)$)

\Rightarrow Existe un T óptimo T_b b y c son hermanos y T_b tiene profundidad máxima

Prueba

Por Lema 1, existe T óptimo $T_b \subset T$ con profundidad máxima

Como T es lleno, c tiene un hermano: x .

Caso 1: $x = b$ Trivial! ✓
Caso 2: $x \neq b$. Aquí tenemos $d_T(b) \leq d_T(x)$ y $f(x) \geq f(b)$

T' swapando x y b . Por el mismo proceso de Lema 3 se demuestra que T' también es óptimo.

Problema: • Σ alfabeto : $|\Sigma| = n$
• $f: \Sigma \rightarrow \mathbb{Z}_0^+$ frecuencia

Sol: Árboles binarios llenos con n hojas, cada una asociada a un carácter. Σ
 $\text{val}(T) = \sum_{x \in \Sigma} f(x) \cdot d_T(x)$

Sol óptima: Aquel T que minimiza $val(T)$.

Elección greedy: Escoger los 2 caracteres con menor frecuencia (b y c)

Subproblema: $\cdot \Sigma^1 = \Sigma - \{b, c\} \cup \{bc\}$ $f_{n-1}(x) = \begin{cases} f_n(b) + f_n(c) & x = bc \\ f_n(x) & x \neq bc \end{cases}$
 $\cdot f_n: \Sigma^1 \rightarrow \mathbb{Z}^+$

elección Greedy:

Existe un T óptimo T_y tiene a b y c en la mayor proporción.

Prueba: Lema 2

Sub problema:

Dado T que tiene a b y c en la mayor proporción.
 T' es el árbol dejado

\Rightarrow T' es óptimo del problema P' .

Por contradicción, sup T' no es óptimo.

$\Rightarrow \exists A'$ óptimo de P' . $val(A') < val(T')$ $\checkmark \odot$

~~3A~~ Construir A :

agregar A' , ubicar la hoja asociada a bc ,
 y le damos 2 hijos: b y c .

Claramente A es sol de P (arb bin llan a loz en Σ)

$$val(A) = val(A') - f_{n-1}(bc) \cdot d_{A'}(bc) + f_n(b) \cdot d_A(b) + f_n(c) \cdot d_A(c)$$

$$val(A) - val(A') = -[f(b) + f(c)] \cdot x + f_n(b) \cdot (x+1) + f_n(c) \cdot (x+1)$$

$$= f_n(b) + f_n(c) = f_n(bc)$$

Similar work

$$\text{val}(T) - \text{val}(T') = f_n(b) + f_n(c)$$

$$\Rightarrow \text{val}(A) - \text{val}(A') = \text{val}(T) - \text{val}(T') -$$

$$- \text{val}(A) = \text{val}(A) - \text{val}(A') + \text{val}(A') \geq \text{val}(T) - \text{val}(T') + \text{val}(A') \Rightarrow \text{val}(A) \geq \text{val}(T)$$