

# Técnicas Fundamentales

## Algoritmos codiciosos II

Segundo semestre 2022

IIC2283

Prof. Nicolás Van Sint Jan

# Recordatorio: Problema de almacenamiento

Sea  $\Sigma$  un alfabeto.

Una  **$\Sigma$ -codificación** es una función  $\tau : \Sigma \rightarrow \{0, 1\}^*$  tal que  $\tau(a) \neq \epsilon$  para todo  $a \in \Sigma$ .

La **extensión**  $\hat{\tau}$  de una  $\Sigma$ -codificación  $\tau$  a todas las palabras  $w \in \Sigma^*$  se define como:

$$\hat{\tau}(w) = \begin{cases} \epsilon & w = \epsilon \\ \tau(a_1) \cdots \tau(a_n) & w = a_1 \cdots a_n \text{ con } n \geq 1 \end{cases}$$

Nos interesa que:

1.  $\hat{\tau}$  sea **inyectiva**.
2.  $\tau$  sea de **largo variable**.

¿Por qué nos interesa esto?

# Recordatorio: Codificación libre de prefijos

## Lema

Si existen  $w_1, w_2 \in \Sigma^*$  tales que  $w_1 \neq w_2$  y  $\hat{\tau}(w_1) = \hat{\tau}(w_2)$ , entonces existen  $a, b \in \Sigma$  tales que  $a \neq b$  y  $\tau(a)$  es un prefijo de  $\tau(b)$

Decimos que una  $\Sigma$ -codificación  $\tau$  es **libre de prefijos** si para cada  $a, b \in \Sigma$  tales que  $a \neq b$  se tiene que  $\tau(a)$  no es un prefijo de  $\tau(b)$ .

## Corolario

Si  $\tau$  es una codificación libre de prefijos, entonces  $\hat{\tau}$  es una **función inyectiva**.

# Recordatorio: Frecuencias relativas de los símbolos

Fije  $w \in \Sigma^*$ .

Para  $a \in \Sigma$  definimos  $\text{fr}_w(a)$  como la **frecuencia relativa** de  $a$  en  $w$ , vale decir

$$\text{fr}_w(a) = \frac{\# \text{ de apariciones de } a \text{ en } w}{|w|}$$

Para una  $\Sigma$ -codificación  $\tau$  definimos el **largo promedio** para  $w$  como:

$$\text{lp}_w(\tau) = \sum_{a \in \Sigma} \text{fr}_w(a) \cdot |\tau(a)|$$

## Problema de optimización a resolver

Dado  $w \in \Sigma^*$ , encontrar una  $\Sigma$ -codificación  $\tau$  libre de prefijos que minimice el valor  $\text{lp}_w(\tau)$

# Recordatorio: Frecuencias relativas de los símbolos

Fije  $w \in \Sigma^*$ .

Para  $a \in \Sigma$  definimos  $\text{fr}_w(a)$  como la **frecuencia relativa** de  $a$  en  $w$ , vale decir

$$\text{fr}_w(a) = \frac{\# \text{ de apariciones de } a \text{ en } w}{|w|}$$

Para una  $\Sigma$ -codificación  $\tau$  definimos el **largo promedio** para  $w$  como:

$$\text{lp}_w(\tau) = \sum_{a \in \Sigma} \text{fr}_w(a) \cdot |\tau(a)|$$

## Problema de optimización a resolver

Dado  $w \in \Sigma^*$ , encontrar una  $\Sigma$ -codificación  $\tau$  libre de prefijos que minimice el valor  $\text{lp}_w(\tau)$

# Recordatorio: Frecuencias relativas de los símbolos

Decimos que  $f : \Sigma \rightarrow (0, 1)$  es una **función de frecuencias relativas** para  $\Sigma$  si se cumple que  $\sum_{a \in \Sigma} f(a) = 1$

Dada una función  $f$  de frecuencias relativas para  $\Sigma$  y una  $\Sigma$ -codificación  $\tau$ , el **largo promedio de  $\tau$  para  $f$**  se define como:

$$\text{lp}_f(\tau) = \sum_{a \in \Sigma} f(a) \cdot |\tau(a)|$$

Problema (correcto) de optimización a resolver

Dado  $f$  una función de frecuencias relativas para  $\Sigma$ , encontrar una  $\Sigma$ -codificación  $\tau$  libre de prefijos que minimice el valor  $\text{lp}_f(\tau)$

# Outline

Algoritmos codiciosos: Codificaciones (cont.)

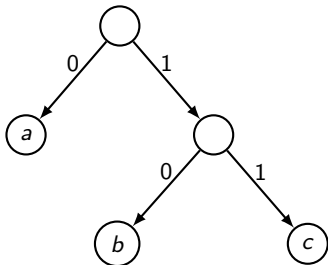
# Outline

Algoritmos codiciosos: Codificaciones (cont.)



# Un ingrediente fundamental: codificaciones como árboles

Representamos la codificación  $\tau(a) = 0$ ,  $\tau(b) = 10$ ,  $\tau(c) = 11$  como un **árbol binario**:



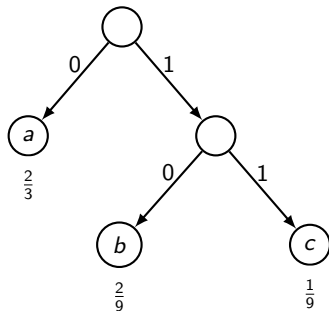
# Un ingrediente fundamental: codificaciones como árboles

Si una  $\Sigma$ -codificación  $\tau$  es **libre de prefijos**, entonces el árbol que la representa satisface las siguientes propiedades.

- Cada hoja tiene como etiqueta un elemento de  $\Sigma$ , y estos son los únicos nodos con etiquetas
- Cada símbolo de  $\Sigma$  es usado exactamente una vez como etiqueta
- Cada arco tiene etiqueta 0 ó 1
- Si una hoja tiene etiqueta  $e$  y las etiquetas de los arcos del camino desde la raíz hasta esta hoja forman una palabra  $w \in \{0, 1\}^*$ , entonces  $\tau(e) = w$

# Codificaciones como árboles y las frecuencias relativas

Podemos agregar al árbol binario que representa una codificación  $\tau$  la información sobre las **frecuencias relativas** dadas por una función  $f$ :



Llamamos a este árbol  $\text{abf}(\tau, f)$ .

# Alguna propiedades importantes

## Ejercicios

Sea  $f$  una función de frecuencias relativas para  $\Sigma$  y  $\tau$  una  $\Sigma$ -codificación libre de prefijos que minimiza la función  $\text{lp}_f(x)$ .

1. Sean  $u$  y  $v$  dos hojas en  $\text{abf}(\tau, f)$  con etiquetas  $a$  y  $b$ , respectivamente. Demuestre que si el camino de la raíz a  $u$  es más corto que el camino de la raíz a  $v$ , entonces  $f(a) \geq f(b)$
2. Demuestre que cada nodo interno en  $\text{abf}(\tau, f)$  tiene dos hijos
3. Sean  $a, b \in \Sigma$  tales que  $a \neq b$ ,  $f(a) \leq f(b)$  y  $f(b) \leq f(e)$  para todo  $e \in (\Sigma \setminus \{a, b\})$ . Demuestre que existe una  $\Sigma$ -codificación  $\tau'$  libre de prefijos tal que  $\text{lp}_f(\tau') = \text{lp}_f(\tau)$  y las hojas con etiquetas  $a$  y  **$b$**  en  $\text{abf}(\tau', f)$  son hermanas
  - Vale decir, existe  $w \in \{0, 1\}^*$  tal que  $\tau'(a) = w0$  y  $\tau'(b) = w1$

# Calculando el mínimo de $lp_f(x)$

Vamos a ver un **algoritmo codicioso** para calcular una  $\Sigma$ -codificación  $\tau$  libre de prefijos que minimiza  $lp_f(x)$

- El algoritmo calcula la **codificación de Huffman**.

El algoritmo tiene los ingredientes mencionados de un algoritmo codicioso.

1. **Función objetivo a minimizar:**  $lp_f(x)$
2. **Función de selección:** elige los dos símbolos de  $\Sigma$  con menor frecuencia relativa, los coloca como hermanos en el árbol binario que representa la  $\Sigma$ -codificación óptima, y continua la construcción con el resto de los símbolos de  $\Sigma$

Será necesario realizar una demostración para probar que el algoritmo es correcto.

# Algoritmo de Huffman

Sea  $f : \Sigma \rightarrow (0, 1)$  una **función de frecuencias relativas**. Supongamos que:

1. Representamos a las funciones como conjuntos de pares ordenados
2.  $f$  tiene al menos dos elementos en el dominio.

El siguiente algoritmo obtiene la  $\Sigma$ -codificación  $\tau$  de Huffman (la que minimiza  $\text{lp}_f(\tau)$ ).

# Algoritmo de Huffman

## **CalcularCodificaciónHuffman( $f$ )**

Sea  $\Sigma$  el dominio de la función  $f$

**if**  $\Sigma = \{a, b\}$  **then return**  $\{(a, 0), (b, 1)\}$

**else**

Sean  $a, b \in \Sigma$  tales que  $a \neq b$ ,  $f(a) \leq f(b)$  y

$f(b) \leq f(e)$  para todo  $e \in (\Sigma \setminus \{a, b\})$

Sea  $c$  un símbolo que **no** aparece en  $\Sigma$

$g := (f \setminus \{(a, f(a)), (b, f(b))\}) \cup \{(c, f(a) + f(b))\}$

$\tau^* := \text{CalcularCodificaciónHuffman}(g)$

$w := \tau^*(c)$

$\tau := (\tau^* \setminus \{(c, w)\}) \cup \{(a, w0), (b, w1)\}$

**return**  $\tau$

# La correctitud de **CalcularCodificaciónHuffman**

## Teorema

Si  $f$  es una función de frecuencias relativas,  $\Sigma$  es el dominio de  $f$  y  $\tau = \mathbf{CalcularCodificaciónHuffman}(f)$ , entonces  $\tau$  es una  $\Sigma$ -codificación libre de prefijos que minimiza la función  $\text{lp}_f(x)$ .

## Demostración

Vamos a realizar la demostración por inducción en  $|\Sigma|$ .

Si  $|\Sigma| = 2$ , entonces la propiedad se cumple trivialmente.

■ ¿Por qué?

Suponga entonces que la propiedad se cumple para un valor  $n \geq 2$ , y suponga que  $|\Sigma| = n + 1$ .



# La demostración del teorema

## Demostración

Sean  $a, b, c, g, \tau^*$  y  $\tau$  definidos como en el código de la llamada **CalcularCodificaciónHuffman**( $f$ ), y sea  $\Gamma$  el dominio de  $g$

Como  $|\Gamma| = n$  y  $\tau^* = \text{CalcularCodificaciónHuffman}(g)$ , por hipótesis de inducción tenemos que  $\tau^*$  es una  $\Gamma$ -codificación libre de prefijos que minimiza la función  $\text{lp}_g(x)$

Dada la definición de  $\tau$  es simple verificar las siguientes propiedades:

- $\tau$  es una codificación libre de prefijos
- Para cada  $e \in (\Sigma \setminus \{a, b\})$  se tiene que  $\tau(e) = \tau^*(e)$
- $|\tau(a)| = |\tau(b)| = |\tau^*(c)| + 1$

# La demostración del teorema

## Demostración

Por contradicción suponga que  $\tau$  NO minimiza el valor de la función  $\text{lp}_f(x)$

- Vale decir, existe una  $\Sigma$ -codificación  $\tau'$  libre de prefijos tal que  $\tau'$  minimiza la función  $\text{lp}_f(x)$  y  $\text{lp}_f(\tau') < \text{lp}_f(\tau)$

Por la definición de  $a$ ,  $b$  y los ejercicios anteriores podemos suponer que existe  $w \in \{0, 1\}^*$  tal que  $\tau'(a) = w \cdot 0$  y  $\tau'(b) = w \cdot 1$ .

- Nótese que  $w \neq \varepsilon$  puesto que  $|\Sigma| \geq 3$ .

A partir de  $\tau'$  defina la siguiente  $\Gamma$ -codificación  $\tau''$ :

$$\tau'' = (\tau' \setminus \{(a, \tau'(a)), (b, \tau'(b))\}) \cup \{(c, w)\}$$

Tenemos que  $\tau''$  es una  $\Gamma$ -codificación libre de prefijos.

- ¿Por qué?

## La relación entre $\text{lp}_g(\tau^*)$ y $\text{lp}_f(\tau)$

$$\begin{aligned}\text{lp}_g(\tau^*) &= \sum_{e \in \Gamma} g(e) \cdot |\tau^*(e)| \\&= \left( \sum_{e \in (\Gamma \setminus \{c\})} g(e) \cdot |\tau^*(e)| \right) + g(c) \cdot |\tau^*(c)| \\&= \left( \sum_{e \in (\Sigma \setminus \{a, b\})} f(e) \cdot |\tau(e)| \right) + (f(a) + f(b)) \cdot |\tau^*(c)| \\&= \left( \sum_{e \in (\Sigma \setminus \{a, b\})} f(e) \cdot |\tau(e)| \right) + f(a) \cdot |\tau^*(c)| + f(b) \cdot |\tau^*(c)| \\&= \left( \sum_{e \in (\Sigma \setminus \{a, b\})} f(e) \cdot |\tau(e)| \right) + f(a) \cdot (|\tau(a)| - 1) + f(b) \cdot (|\tau(b)| - 1) \\&= \left( \sum_{e \in \Sigma} f(e) \cdot |\tau(e)| \right) - (f(a) + f(b)) \\&= \text{lp}_f(\tau) - (f(a) + f(b))\end{aligned}$$

## La relación entre $\text{lp}_g(\tau'')$ y $\text{lp}_f(\tau')$

$$\begin{aligned}\text{lp}_g(\tau'') &= \sum_{e \in \Gamma} g(e) \cdot |\tau''(e)| \\&= \left( \sum_{e \in (\Gamma \setminus \{c\})} g(e) \cdot |\tau''(e)| \right) + g(c) \cdot |\tau''(c)| \\&= \left( \sum_{e \in (\Sigma \setminus \{a, b\})} f(e) \cdot |\tau'(e)| \right) + (f(a) + f(b)) \cdot |\tau''(c)| \\&= \left( \sum_{e \in (\Sigma \setminus \{a, b\})} f(e) \cdot |\tau'(e)| \right) + f(a) \cdot |\tau''(c)| + f(b) \cdot |\tau''(c)| \\&= \left( \sum_{e \in (\Sigma \setminus \{a, b\})} f(e) \cdot |\tau'(e)| \right) + f(a) \cdot (|\tau'(a)| - 1) + f(b) \cdot (|\tau'(b)| - 1) \\&= \left( \sum_{e \in \Sigma} f(e) \cdot |\tau'(e)| \right) - (f(a) + f(b)) \\&= \text{lp}_f(\tau') - (f(a) + f(b))\end{aligned}$$

# Obteniendo una contradicción

## Demostración

Tenemos entonces que

$$\begin{aligned} \text{lp}_g(\tau'') &= \text{lp}_f(\tau') - (f(a) + f(b)) \\ &< \text{lp}_f(\tau) - (f(a) + f(b)) \\ &= \text{lp}_g(\tau^*) \end{aligned}$$

Concluimos entonces que  $\tau^*$  NO minimiza la función  $\text{lp}_g(x)$ , lo cual contradice la hipótesis de inducción. □

## Dos comentarios finales

La siguiente función calcula la **codificación de Huffman** teniendo como entrada una palabra  $w$ :

**CalcularCodificaciónHuffman**( $w$ )

**if**  $w = \varepsilon$  **then return**  $\emptyset$

**else**

$\Sigma :=$  conjunto de símbolos mencionados en  $w$

**if**  $\Sigma = \{a\}$  **then return**  $\{(a, 0)\}$

**else return** **CalcularCodificaciónHuffman**( $\text{fr}_w$ )

¿Cuál es la complejidad de **CalcularCodificaciónHuffman**( $f$ )?

**R:**  $\mathcal{O}(n \log n)$  (considerando que  $|f| \in \Theta(|\Sigma|)$ ).