

TEORÍA DE ALGORITMOS  
(75.29) CURSO BUCHWALD - GENENDER

# Trabajo Práctico 2

## Que parezca programación dinámica

5 de mayo de 2025

Julen Leonel Gaumard  
111379

Noelia Salvatierra  
100116

Franco Macke  
105974

## 1. Análisis del problema

Tenemos:

- Una cadena encriptada  $S$  sin separaciones
- Un diccionario de palabras válidas en español

Debemos determinar:

- Si es posible separar  $S$  en una secuencia de palabras de  $D$ . En caso afirmativo, cuál sería esa separación.
- Si podemos determinar que los primeros  $i$  caracteres forman una secuencia válida de palabras, entonces solo necesitamos verificar si el resto de la cadena también puede separarse en palabras válidas.
- Para cada posición en la cadena, podemos probar todas las posibles palabras que podrían comenzar desde allí y verificar si el resto de la cadena puede separarse recursivamente.

### Ecuación de recurrencia

Sea `oracion[0..n-1]` la cadena a segmentar. Definimos `existencia_parcial` tal que `existencia_parcial[i]` es verdadero si la subcadena de longitud  $i$  de la cadena (es decir, `oracion[0..i-1]`) puede segmentarse completamente en palabras del diccionario.

Entonces, definimos la ecuación de recurrencia tal que:

$$\text{existencia\_parcial}[i] = \bigvee_{j=0}^{i-1} (\text{existencia\_parcial}[j] \wedge \text{oracion}[j : i] \in \text{diccionario})$$

Esto significa que `existencia_parcial[i]` va a ser verdadero si existe al menos un índice  $j < i$  tal que:

- el prefijo `oracion[0..j-1]` puede segmentarse (es decir, `existencia_parcial[j] = True`)
- la subcadena `oracion[j..i-1]` pertenece al diccionario

El caso base es:

$$\text{existencia\_parcial}[0] = \text{True}$$

ya que una cadena vacía la consideramos una palabra válida.

## 2. Demostración

La ecuación de recurrencia obtenida es:

$$\text{existencia\_parcial}[i] = \bigvee_{j=0}^{i-1} (\text{existencia\_parcial}[j] \wedge \text{oracion}[j : i] \in \text{diccionario})$$

Donde:

$\text{existencia\_parcial}[i]$  es una tabla booleana que va a ser una  $P(n)$  proposición, con  $i \in [0, n]$  tal que:

$\text{existencia\_parcial}[i] = \text{True} \iff$  la subcadena  $\text{oracion}[0 : i]$  puede segmentarse completamente en palabras del diccionario.

Demostraremos que esta ecuación llega a la solución óptima, por el principio de inducción matemática y usando el método directo:

**Paso base:**  $i = 0$

La subcadena vacía  $\text{oracion}[0 : 0]$  se puede segmentar (no contiene caracteres), por lo tanto:

$$\text{existencia\_parcial}[0] = \text{True}$$

**Hipótesis verdadera:** Supongamos que para todo  $k < i$ ,  $\text{existencia\_parcial}[k] = \text{True}$  si y solo si el prefijo  $\text{oracion}[0 : k]$  PUEDE SEGMENTARSE completamente en palabras del diccionario.

**Paso inductivo:** Demostraremos que bajo esta hipótesis,  $\text{existencia\_parcial}[i]$  también se establece correctamente.

Según la ecuación de recurrencia:

$$\text{existencia\_parcial}[i] = \bigvee_{j=0}^{i-1} (\text{existencia\_parcial}[j] \wedge \text{oracion}[j : i] \in D)$$

Esto significa que:

- Existe un índice  $j < i$  tal que  $\text{existencia\_parcial}[j] = \text{True}$  (es decir,  $\text{oracion}[0 : j]$  puede segmentarse), y
- La subcadena  $\text{oracion}[j : i]$  pertenece al diccionario.

Entonces, se puede construir la segmentación de  $\text{oracion}[0 : i]$  como la concatenación de dos partes válidas:

$$\text{oracion}[0 : j] + \text{oracion}[j : i]$$

y por lo tanto,  $\text{existencia\_parcial}[i] = \text{True}$ .

Por otro lado, si no existe ningún  $j < i$  que cumpla esa condición, entonces no existe una forma de segmentar  $\text{oracion}[0 : i]$ , y correctamente:

$$\text{existencia\_parcial}[i] = \text{False}$$

**Conclusión:** Por el principio de inducción matemática, la proposición  $P(n)$  (tabla  $\text{existencia\_parcial}$  está correctamente construida para todo  $i \in [0, n]$ ).

### 3. Complejidad teorica

Primero, vamos a analizar el código analíticamente:

```
1 def procesar_texto(oraciones, diccionario):
2     """
3     Procesa una lista de oraciones y un diccionario, y devuelve una lista de
4     oraciones segmentadas.
5     :param oraciones: Lista de oraciones a procesar.
6     :param diccionario: Diccionario utilizado para la segmentación.
7     """
8     resultado = []
9
10    for oracion in oraciones:
11        oracion_segmentada = segmentar_oracion(oracion, diccionario)
12
13        if len(" ".join(oracion_segmentada)) < len(oracion):
14            resultado.append("No es un mensaje")
15        else:
16            resultado.append(" ".join(oracion_segmentada))
17
18    return resultado
19
20 def segmentar_oracion(oracion, diccionario):
21     n = len(oracion)
22     conjunto_diccionario = set(diccionario)
23
24     existencia_parcial = [False] * (n + 1)
25     existencia_parcial[0] = True
26
27     path = [None] * (n + 1)
28
29     for i in range(1, n + 1):
30         for j in range(i):
31             if existencia_parcial[j] and oracion[j:i] in conjunto_diccionario:
32                 existencia_parcial[i] = True
33                 path[i] = (j, oracion[j:i])
34                 break
35
36     if not existencia_parcial[n]:
37         return []
38
39     return reconstruir_segmentacion(path, n)
40
41
42 def reconstruir_segmentacion(path, n):
43     """
44     Reconstruye la segmentación de una oración a partir del camino dado.
45     :param path: Lista de tuplas que representan el camino de segmentación.
46     :param n: Longitud de la oración original.
47     :return: Lista de palabras segmentadas.
48     """
49     resultado = []
50     idx = n
51     while idx > 0:
52         j, palabra = path[idx]
53         resultado.append(palabra)
54         idx = j
55
56     return resultado[::-1]
```

Listing 1: Algoritmo principal

Analizando línea por línea:

- `resultado = []`: asignación Complejidad:  $\mathcal{O}(1)$ .
- Bucle principal `for`:

- Se ejecuta  $m$  veces (una por cada oracion), siendo  $m$  el tamaño de oraciones.
- En cada iteración:
  - `segmentar_oracion(...)`: (se justifica con detalle más adelante).  $\mathcal{O}(k + n^3)$  con  $k$ , el tamaño del diccionario.
  - `condicional if`:  $\mathcal{O}(1)$ . (los condicionales se ejecutan a tiempo constante)
  - `append(...)`: tiempo constante.  $\mathcal{O}(1)$ .
  - `condicional else`:  $\mathcal{O}(1)$ .
  - `append(...)`:  $\mathcal{O}(1)$ .
- Total por iteración:  $\mathcal{O}(k + n^3)$ .
- Complejidad total del bucle:  $\mathcal{O}(m * (k + n^3))$ .
- `return resultado`:  $\mathcal{O}(1)$ .

#### Complejidad total del algoritmo:

$$\mathcal{O}(1) + \mathcal{O}(m * (k + n^3)) + \mathcal{O}(1) = \boxed{\mathcal{O}(m * (k + n^3))}$$

Entonces, analíticamente podemos decir que el algoritmo es  $\mathcal{O}(m * (k + n^3))$  con  $m = \text{len}(\text{oraciones})$ ,  $n = \text{len}(\text{oracion})$  y  $k = \text{len}(\text{diccionario})$ .

Observamos que si el diccionario de palabras supera a  $n$ , la complejidad se puede aproximar a:  $\mathcal{O}(m * k)$ .

Si en cambio,  $m$  es mucho mas grande que el diccionario de palabras, la complejidad se puede aproximar a:  $\mathcal{O}(m)$ , donde  $m$  es la cantidad de oraciones a segmentar.

`Segmentar_oracion(...)`:

- `len(oracion)`: Complejidad:  $\mathcal{O}(n)$ .
- `set(diccionario)`: Complejidad:  $\mathcal{O}(k)$  con  $k$  el tamaño del diccionario.
- `existencia_parcial = [False] * (n + 1)`: Complejidad:  $\mathcal{O}(n)$ .
- `existencia_parcial[0] = True`: Complejidad:  $\mathcal{O}(1)$ .
- `path = [None] * (n + 1)`: Complejidad:  $\mathcal{O}(n)$ .
- Bucle principal `for`:
  - Se ejecuta  $n$  veces (una por  $i$  caracter), siendo  $n$  el tamaño de oración.
  - En cada iteración:
    - `bucle secundario for`:
      - ◊ Se ejecuta  $i$  veces (una por  $j$  caracter)
      - ◊ En cada iteración:
        - ◊ `condicional if con slice oracion[j:i]`: Los slices no son de tiempo constante. Complejidad  $\mathcal{O}(i)$ .
        - ◊ `resto de asignaciones`: Complejidad  $\mathcal{O}(1)$ .
        - ◊ `break`: Complejidad  $\mathcal{O}(1)$ .
        - ◊ Total por iteración:  $\mathcal{O}(i)$ .
      - ◊ Complejidad total del bucle secundario:  $\mathcal{O}(i^2)$ .
    - Total por iteración:  $\mathcal{O}(i^2)$ .
  - Complejidad total del bucle principal:  $\mathcal{O}((i^2) * n)$ . Ahora, en el peor caso,  $i = n$  (sino llegamos al `break`), entonces en este caso, podemos aproximar la complejidad a  $\mathcal{O}(n^3)$ .

- `return reconstruir_segmentacion(...)`: función que reconstruye el mensaje con sentido”, dada la solución por programación dinámica y tiene asignaciones, llamadas a `append(...)`, que son de tiempo constante, y un bucle `while` que en el peor caso, se recorre hasta que `idx` menor o igual a 0, con `idx = n`, entonces la complejidad es:  $\mathcal{O}(n)$ .

**Complejidad total segmentar oracion(...):**

$$\mathcal{O}(n) + \mathcal{O}(k) + \mathcal{O}(n^3) + \mathcal{O}(n) = \boxed{\mathcal{O}(k + n^3)}$$

## 4. Ejemplos de ejecución

**Oraciones generadas:** zunaoirá

**Resultados del procesamiento:**

Evaluando subcadena: z

Es válido hasta la posición 0?: True

Existe en el diccionario?: False

Evaluando subcadena: zu

Es válido hasta la posición 0?: True

Existe en el diccionario?: False

Evaluando subcadena: u

Es válido hasta la posición 1?: False

Existe en el diccionario?: False

Evaluando subcadena: zun

Es válido hasta la posición 0?: True

Existe en el diccionario?: False

Evaluando subcadena: un

Es válido hasta la posición 1?: False

Existe en el diccionario?: False

Evaluando subcadena: n

Es válido hasta la posición 2?: False

Existe en el diccionario?: False

Evaluando subcadena: zuna

Es válido hasta la posición 0?: True

Existe en el diccionario?: True

Evaluando subcadena: zunao

Es válido hasta la posición 0?: True

Existe en el diccionario?: False

Evaluando subcadena: unao

Es válido hasta la posición 1?: False

Existe en el diccionario?: False

Evaluando subcadena: nao

Es válido hasta la posición 2?: False

Existe en el diccionario?: False

Evaluando subcadena: ao

Es válido hasta la posición 3?: False

Existe en el diccionario?: False

Evaluando subcadena: o

Es válido hasta la posición 4?: True

Existe en el diccionario?: False

Evaluando subcadena: zunaoi

Es válido hasta la posición 0?: True

Existe en el diccionario?: False

Evaluando subcadena: unaoi

Es válido hasta la posición 1?: False

Existe en el diccionario?: False

Evaluando subcadena: naoi

Es válido hasta la posición 2?: False

Existe en el diccionario?: False

Evaluando subcadena: ao i

Es válido hasta la posición 3?: False

Existe en el diccionario?: False

Evaluando subcadena: oi

Es válido hasta la posición 4?: True

Existe en el diccionario?: False

Evaluando subcadena: i

Es válido hasta la posición 5?: False

Existe en el diccionario?: False

Evaluando subcadena: zunaoir

Es válido hasta la posición 0?: True

Existe en el diccionario?: False

Evaluando subcadena: unaoir

Es válido hasta la posición 1?: False

Existe en el diccionario?: False

Evaluando subcadena: naoir

Es válido hasta la posición 2?: False

Existe en el diccionario?: False

Evaluando subcadena: aoir

Es válido hasta la posición 3?: False

Existe en el diccionario?: False

Evaluando subcadena: oir

Es válido hasta la posición 4?: True

Existe en el diccionario?: False

Evaluando subcadena: ir

Es válido hasta la posición 5?: False

Existe en el diccionario?: False

Evaluando subcadena: r

Es válido hasta la posición 6?: False

Existe en el diccionario?: False

Evaluando subcadena: zunaoirá

Es válido hasta la posición 0?: True

Existe en el diccionario?: False

Evaluando subcadena: unaoirá  
Es válido hasta la posición 1?: False  
Existe en el diccionario?: False

Evaluando subcadena: naoirá  
Es válido hasta la posición 2?: False  
Existe en el diccionario?: False

Evaluando subcadena: aoirá  
Es válido hasta la posición 3?: False  
Existe en el diccionario?: False

Evaluando subcadena: oirá  
Es válido hasta la posición 4?: True  
Existe en el diccionario?: True

Oración: zunaoirá    Resultado: zuna oirá

## 5. Análisis de la complejidad por cuadrados minimos

Vamos a corroborar la complejidad calculada anteriormente de manera empírica. Para ello, seteamos la cantidad de palabras y el tamaño de las oraciones a descifrar relativamente grandes. Los casos generados se hacen utilizando palabras del diccionario generadas aleatoriamente de esta manera:

```
1 def generate_random_dictionary(size, min_length=3, max_length=MAX_WORD_LENGTH):
2     """Genera un diccionario aleatorio de palabras."""
3     dictionary = []
4     for _ in range(size):
5         word_length = random.randint(min_length, max_length)
6         word = ''.join(random.choice(string.ascii_lowercase) for _ in range(
7             word_length))
8         dictionary.append(word)
9     return dictionary
10
11 def generate_valid_text(dictionary, length):
12     """Genera un texto aleatorio usando palabras del diccionario."""
13     return ''.join(random.choice(dictionary) for _ in range(length))
```

Listing 2: Generación de diccionario y texto válido aleatorios

Primero, queremos comprobar que efectivamente, la complejidad de `segmentar_oracion()` es  $\mathcal{O}(n^3 + k)$ .

Hacemos el análisis de cantidad de palabras del diccionario  $k$  constante.

### Constantes de longitud de palabras

`MAX_WORD_LENGTH = 400` (Longitud máxima de las palabras)

`NUMBER_OF_WORDS_IN_ONE_SENTENCE = 100` (Cantidad de palabras por oracion)

### Cantidad de oraciones por iteración

`sizes = [50, 100, 125, 150, 170, 200, 250, 300, 400]`

`NUMBER_OF_SENTENCES = 1` (Cantidad de oraciones por caso, m en nuestro análisis)



## Gráfico de complejidad temporal

A continuación se muestra el gráfico de complejidad temporal empírica usando las variables anteriormente mencionadas, acompañado de la curva que mejor ajusta la tendencia observada.

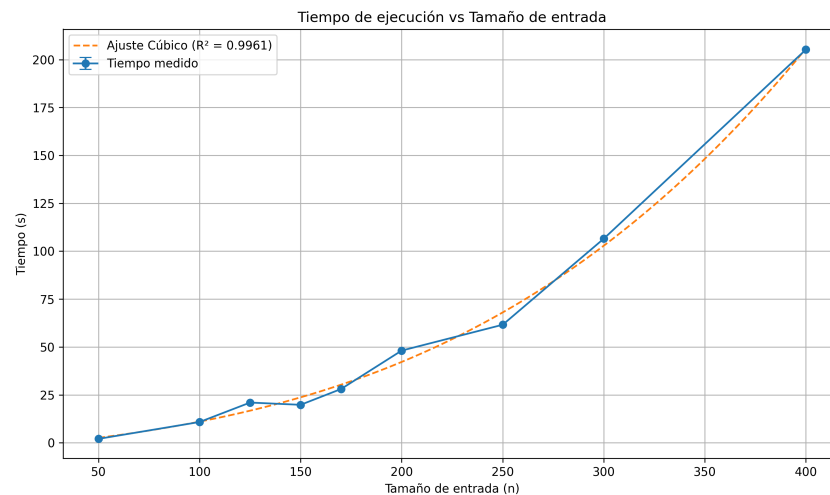


Figura 1: Ajuste de complejidad temporal con curva cúbica

## Ajuste por mínimos cuadrados

La curva ajustada corresponde a la siguiente expresión:

$$y = 2,11 \times 10^{-6} \cdot x^3 + 2,12 \times 10^{-4} \cdot x^2 + 1,00 \cdot x - 3,23$$

Efectivamente, observamos que la complejidad de `segmentar_oracion()` resulta  $\mathcal{O}(n^3)$ .

Subiendo  $m$  (`NUMBER_OF_SENTENCES`) a un valor proporcional a  $n$ , aumentaría la complejidad a:

$$T(n, m) = m \cdot \mathcal{O}(n^3)$$

Si  $m = \mathcal{O}(n)$ , entonces:

$$T(n, m) = \mathcal{O}(n) \cdot \mathcal{O}(n^3) = \mathcal{O}(n^4)$$

## 6. Variabilidad de los valores

### Complejidad con diccionario

Ahora, mantenemos  $n$  constante y variamos la cantidad de palabras del diccionario.

```
1 DICTIONARY_SIZES = [1000, 2000, 5000, 7000, 10000, 15000, 20000]
```

Listing 3: Tamaños del diccionario utilizados

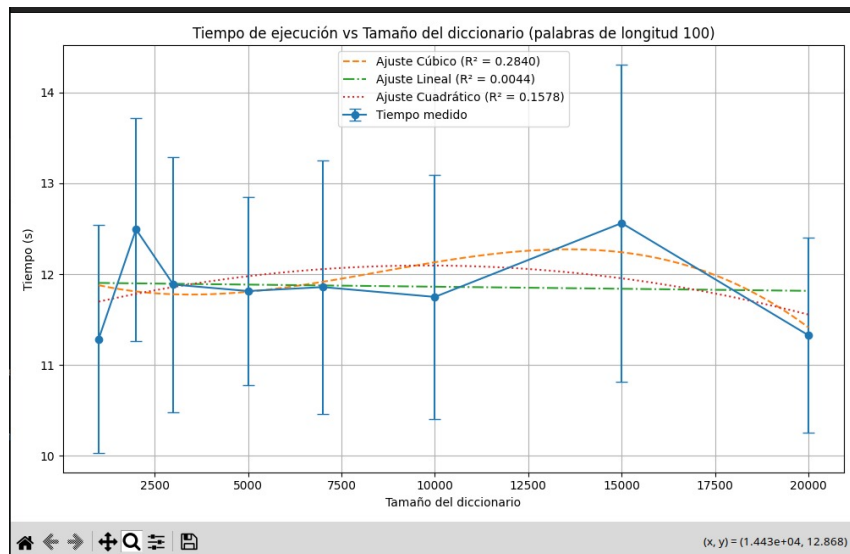


Figura 2: Mantengo  $n$  constante y aumento  $k$

Podemos ver que tarda prácticamente lo mismo con un diccionario de longitud 1.000 y uno de 20.000. Podemos ver que cuando  $n$  es grande,  $k$  termina siendo absorbido en la complejidad.

### Variante $n$ chico y $k$ muy grande

Para mostrar la dependencia con la cantidad de palabras del diccionario, proponemos el caso donde  $n$  es chico y  $k$  muy grande. Esto resulta en el gráfico siguiente, siendo la mejor aproximación algo lineal.

```
1 HUGE_DICTIONARY_SIZES = [num for num in range(10_000, 1_000_000, 50_000)]
```

Listing 4: Generación de tamaños grandes de diccionario

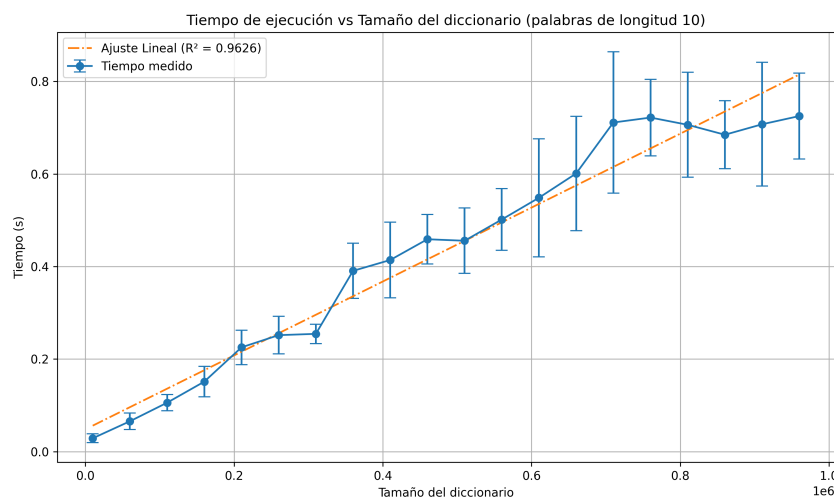


Figura 3: Mantengo  $n$  chica y constante aumentando  $k$

Claramente,  $k$  no puede ser obviada, ya que si es mucho más grande que  $n$ , afecta a la complejidad temporal.

## Variaciones de k

Cuando  $k \propto n^3$ :

$$O(k + n^3) = O(c \cdot n^3 + n^3) = O(n^3 + n^3) = O(2n^3) = O(n^3)$$

$$\therefore O(k + n^3) = O(n^3)$$

Cuando  $k \ll n^3$ :

$$O(k + n^3) = O(n^3)$$

El término  $k$  es despreciable frente a  $n^3$

Cuando  $k \gg n^3$ :

$$O(k + n^3) = O(k)$$

El término  $n^3$  es despreciable frente a  $k$