

REPORTE TAREA 2 y 3

ALGORITMOS Y COMPLEJIDAD

«Explorando la Distancia entre Cadenas, una Operación a la Vez»

Joaquín Domínguez Bustos

17 de noviembre de 2024

23:47

Resumen

Este informe aborda el problema de la distancia de Damerau-Levenshtein mediante dos enfoques principales: algoritmos de fuerza bruta y programación dinámica. Se analiza cómo los costos asociados a operaciones como eliminación, inserción, reemplazo y transposición afectan la complejidad temporal de las soluciones propuestas. El trabajo incluye un diseño detallado de los algoritmos, su implementación en código y la evaluación experimental en diversos conjuntos de datos. Los resultados muestran una clara ventaja en eficiencia para la programación dinámica, aunque a costa de un mayor uso de memoria, mientras que la fuerza bruta puede ser útil en casos específicos con cadenas vacías.

Índice

1. Introducción	2
2. Diseño y Análisis de Algoritmos	3
3. Implementaciones	6
4. Experimentos	7
5. Conclusiones	12
6. Condiciones de entrega	13
A. Apéndice 1	14

1. Introducción

En el ámbito de las Ciencias de la Computación, el análisis y diseño de algoritmos constituye un área fundamental, ya que establece las bases para resolver problemas complejos de manera eficiente. Uno de los problemas clásicos en este campo es la comparación de cadenas, esencial en aplicaciones como la corrección ortográfica, el reconocimiento de patrones y la bioinformática. La distancia de Damerau-Levenshtein, que mide el costo mínimo para transformar una cadena en otra mediante operaciones de inserción, eliminación, reemplazo y transposición, proporciona una métrica versátil y ampliamente utilizada.

A pesar de su relevancia, este problema plantea desafíos significativos en términos de eficiencia computacional, especialmente cuando se utilizan algoritmos de alta complejidad como la fuerza bruta. Las soluciones modernas, como las basadas en programación dinámica, prometen una mejora considerable, aunque con un costo adicional de memoria. Sin embargo, existe una brecha en la comprensión detallada del impacto de las operaciones individuales y los escenarios específicos sobre el rendimiento de estos algoritmos.

El propósito de este informe es estudiar y comparar ambos enfoques, fuerza bruta y programación dinámica, en el contexto del cálculo de la distancia de Damerau-Levenshtein. A través de un diseño detallado, implementaciones prácticas y experimentos controlados, se busca determinar cómo los costos asociados a cada operación afectan la complejidad temporal y espacial de las soluciones. Además, se evalúan casos particulares, como cadenas vacías o con caracteres repetidos, para identificar patrones de rendimiento y posibles áreas de optimización.

2. Diseño y Análisis de Algoritmos

- La solución para ambos algoritmos ([algoritmo 1](#) y [algoritmo 2](#)) está basada en el uso de índices (i, j) con los cuales se recorren las cadenas $S1$ y $S2$ respectivamente para realizar el calculo de todas las posibles operaciones con las cuales se puede transformar la cadena $S1$ en la cadena $S2$ (eliminación, inserción, reemplazo y transposición).

2.1. Fuerza Bruta

Algoritmo 1: Algoritmo de Fuerza Bruta para la Transformación de Cadenas

```

1 Procedure BF( $S1, S2, i, j$ )
2   if  $i \geq longitud(S1)$  and  $j \geq longitud(S2)$  then
3     return 0
4   if  $i \geq longitud(S1)$  then
5     return COSTO_INS( $S2[j]$ ) + BF( $S1, S2, i, j+1$ )
6   if  $j \geq longitud(S2)$  then
7     return COSTO_DEL( $S1[i]$ ) + BF( $S1, S2, i+1, j$ )
8    $elim \leftarrow$  COSTO_DEL( $S1[i]$ ) + BF( $S1, S2, i+1, j$ )
9    $ins \leftarrow$  COSTO_INS( $S2[j]$ ) + BF( $S1, S2, i, j+1$ )
10   $remp \leftarrow$  COSTO_SUB( $S1[i], S2[j]$ ) + BF( $S1, S2, i+1, j+1$ )
11   $transp \leftarrow \infty$ 
12  if  $i+1 < longitud(S1)$  and  $j+1 < longitud(S2)$  and  $S1[i] = S2[j+1]$  and  $S1[i+1] = S2[j]$  then
13     $transp \leftarrow$  COSTO_TRANS( $S1[i], S1[i+1]$ ) + BF( $S1, S2, i+2, j+2$ )
14  return  $\min\{elim, ins, remp, transp\}$ 

```

2.1.1. Ejemplo de ejecución

Como ejemplo tomemos las cadenas $S1$: 'abc' $S2$: 'abs' El algoritmo reemplazaria 'a' por 'a', luego reemplazaria 'b' por 'b' y posteriormente eliminaria 'c' e insertaria 's'

2.1.2. Análisis de complejidad espacial y temporal

En cada llamada que se realiza se pueden generar hasta 4 subproblemas nuevos, así tambien es necesario recorrer cada cadena hasta el final de la misma mediante el uso de los índices (i, j) , los cuales llegaran a un valor final (n, m) en donde n y m son los largos de $S1$ y $S2$ respectivamente, por lo que el número total de operaciones será aproximadamente $4^{\min(n, m)}$, dando así que la complejidad temporal pertenece a $O(4^{\min(n, m)})$. Por otra parte, la complejidad espacial asociada a esta implementacion sería $O(n + m)$ puesto que no se utiliza ninguna estructura de datos adicional al tamaño de la entrada y para las llamadas recursivas se apoya en el stack.

2.2. Programación Dinámica

2.2.1. Descripción de la solución recursiva

Por cada llamada a la función es necesario evaluar los cuatro casos posibles (eliminación, inserción, reemplazo y transposición) y escoger el de menor costo, si guardamos el costo asociado a esa llamada este nos puede servir para evitar realizar el mismo cálculo si se da el caso, de esta manera se asegura que cada llamada que posea los mismos parámetros se resuelva una única vez, esto implica que mediante la consulta recurrente de valores podemos disminuir el costo temporal del problema en comparación a la utilización de fuerza bruta y obtendremos exactamente la misma respuesta final.

2.2.2. Relación de recurrencia

- *(condicion)* Si el valor de $Cache[i][j]$ es distinto de -1 se retorna el valor de $Cache[i][j]$
- *(caso base)* Si el valor de i es igual al largo de $S1$ y el valor de j igual al largo de $S2$ entonces se retorna 0
- *(condicion)* Si el valor de i iguala el largo de $S1$ entonces solo se realizan inserciones por cada carácter faltante por cubrir de $S2$, devolviendo el costo de estas operaciones.
- *(condicion)* Si el valor de j iguala el largo de $S2$ entonces solo se realizan eliminaciones por cada carácter sobrante de $S1$, devolviendo el costo de estas operaciones.
- *(Recurrencia)*

$$DP(S1, S2, i, j) = Cache[i][j] = \min \left\{ \begin{array}{l} \text{costo_del}(S1[i]) + DP(S1, S2, i + 1, j), \\ \text{costo_ins}(S2[j]) + DP(S1, S2, i, j + 1), \\ \text{costo_sub}(S1[i], S2[j]) + DP(S1, S2, i + 1, j + 1), \\ \text{costo_trans}(S1[i], S1[i + 1]) + DP(S1, S2, i + 2, j + 2) \end{array} \right\}$$

2.2.3. Identificación de subproblemas

Los subproblemas son las combinaciones de las posiciones i y j que representan el costo mínimo para transformar la subcadena de $S1$ desde i a la subcadena de $S2$ desde j .

2.2.4. Estructura de datos y orden de cálculo

Se utiliza el enfoque top-down una tabla de $Cache$ para almacenar los resultados, que son calculados recursivamente, garantizando que cada subproblema solo se evalúe una vez.

2.2.5. Algoritmo utilizando programación dinámica

Algoritmo 2: Algoritmo de Programación Dinámica para la Transformación de Cadenas

Input: Cadenas $S1$, $S2$, índices i , j , y matriz de cache de tamaño $(n+1) \times (m+1)$ inicializada en -1

```

1 Procedure DP( $S1, S2, i, j, Cache$ )
2   if  $Cache[i][j] \neq -1$  then
3     return  $Cache[i][j]$ 
4   if  $i = \text{largo}(S1)$  and  $j = \text{largo}(S2)$  then
5     return  $Cache[i][j] = 0$ 
6   if  $i = \text{largo}(S1)$  then
7     return  $Cache[i][j] = \text{costo\_ins}(S2[j]) + DP(S1, S2, i, j + 1, Cache)$ 
8   if  $j = \text{largo}(S2)$  then
9     return  $Cache[i][j] = \text{costo\_del}(S1[i]) + DP(S1, S2, i + 1, j, Cache)$ 
10   $elim \leftarrow \text{costo\_del}(S1[i]) + DP(S1, S2, i + 1, j, Cache);$ 
11   $ins \leftarrow \text{costo\_ins}(S2[j]) + DP(S1, S2, i, j + 1, Cache);$ 
12   $remp \leftarrow \text{costo\_sub}(S1[i], S2[j]) + DP(S1, S2, i + 1, j + 1, Cache);$ 
13   $transp \leftarrow \infty;$ 
14  if  $i + 1 < \text{largo}(S1)$  and  $j + 1 < \text{largo}(S2)$  and  $S1[i] = S2[j + 1]$  and  $S1[i + 1] = S2[j]$  then
15     $transp \leftarrow \text{costo\_trans}(S1[i], S1[i + 1]) + DP(S1, S2, i + 2, j + 2, Cache)$ 
16  return  $Cache[i][j] = \min(elim, ins, remp, transp)$ 

```

2.2.6. Ejemplo de ejecución

Como ejemplo tomemos las cadenas $S1$: 'abc' $S2$: 'abs' El algoritmo reemplazaria 'a' por 'a', luego reemplazaria 'b' por 'b' y posteriormente eliminaria 'c' e insertaria 's'

2.2.7. Análisis de complejidad espacial y temporal

La implementación de programación dinámica utiliza una tabla de caché para almacenar el costo mínimo de edición para cada par (i, j) , permitiendo así que cada estado (eliminación, inserción, sustitución y transposición) realice un número constante de operaciones implicando un costo temporal $O(1)$ ya que solo consideramos la consulta al costo, aunque, en el peor caso, será necesario llenar por completo la tabla de caché, la cual tiene dimensiones $(n + 1) \times (m + 1)$, donde n y m son los largos de $S1$ y $S2$, respectivamente. Esto significa que hay $O(n \times m)$ estados posibles, y cada uno se evaluará a lo sumo una vez, por lo que la complejidad temporal es $O(n \times m)$. En cuanto al análisis espacial, el uso de la tabla de caché la que utiliza mayor espacio, lo que implica que la complejidad espacial también es $O(n \times m)$.

3. Implementaciones

A continuacion se detallan los archivos utilizados y el contenido de cada uno, además en el apendice se puede encontrar el link al repositorio de github.

- **Archivos .cpp y .hpp**

- **DP.cpp:** Algoritmo de programacion dinamica.
- **BF.cpp:** Algoritmo de fuerza bruta.
- **utils.hpp:** Funcion de lectura de archivos para los costos, funciones de costos, funcion de tamaño espacial, funcion de reconstruccion.

- **Generadores**

- **GeneradorData.py:** Genera los datasets a utilizar.
- **GenerarCostos.py:** Genera los costos de las operaciones de insercion, eliminacion, reemplazo y transposicion.
- **Graficos.py:** Genera los graficos con los resultados almacenados en BF.txt y DP.txt.

- **Archivos .txt**

- **cost_operacion.txt:** estos archivos contienen las matrices de costos generadas.
- **DP.txt y BF.txt:** contienen los datos de salida de cada ejecucion
- **(datasets).txt:** estos archvos contienen los casos de prueba utilizados.

4. Experimentos

Todos los experimentos se desarrollaron en el siguiente ambiente computacional:

■ Hardware

- **Procesador:** AMD Ryzen 5 4600H 3.0GHz
- **RAM:** 16GB Dual Channel 3200 MT/s DDR4 SODIMM
- **Almacenamiento:** 500GB SSDNVMe

■ Software

- **SO:** Windows 11 v10.0.22631.4460
- **Entorno:** WSL v2.3.26.0 con Ubuntu 14.2.0 y g++ 14.2.0
- **Python:** v3.12.3 con librerías pandas, matplotlib.pyplot, numpy, pathlib, random.

4.1. Dataset (casos de prueba)

■ 1. IgualLargo.txt

- Este archivo contiene pares de cadenas de igual longitud. Cada par está identificado por un número que indica la longitud de las cadenas, seguido de las cadenas mismas.
- La utilización de este dataset se utiliza a modo de estudiar como varía el tiempo de ejecución de ambos algoritmos sobre casos aleatorios.

■ 2. LetrasRepetidas.txt

- Este archivo incluye pares de cadenas donde los caracteres de cada cadena son iguales y repetidos.
- Este dataset se utiliza para estudiar el peor caso posible teóricamente, en el cual se pueden realizar todas las operaciones en cada llamada

■ 3. Transposiciones.txt

- Este archivo contiene pares de cadenas que difieren únicamente en una o más transposiciones de caracteres consecutivos.
- Este dataset se utiliza a modo de estudiar la complejidad adicional que proporcionan las transposiciones y la optimización del costo.

■ 4. VacioS1.txt y VacioS2.txt

- Estos archivos incluyen cadenas donde una de las cadenas está vacía, mientras que la otra tiene diferentes longitudes.
- Permite estudiar casos donde todas las operaciones son solamente inserciones o eliminaciones.

4.2. Resultados

A continuacion se discuten los resultados obtenidos por cada dataset

Análisis de Rendimiento - Dataset: IgualLargo

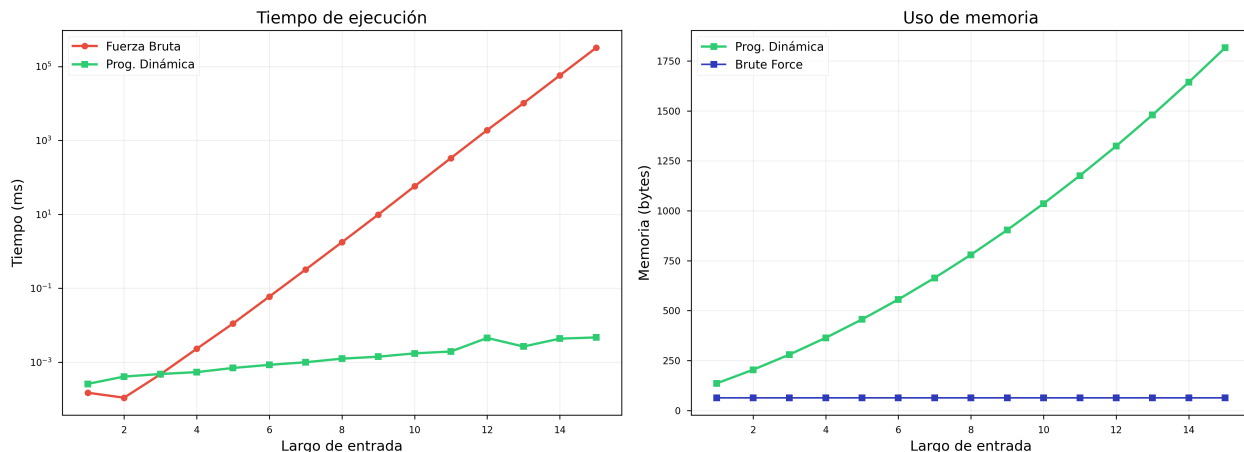


Figura 1: Cadenas aleatorias de largo creciente

En el grafico de la izquierda se puede observar la inmensa diferencia en tiempo de ejecucion de ambos algoritmos a medida que el tamaño de la entrada crece, aunque la menor complejidad temporal que maneja programacion dinamica requiere una mayor complejidad espacial.

Análisis de Rendimiento - Dataset: Transposiciones

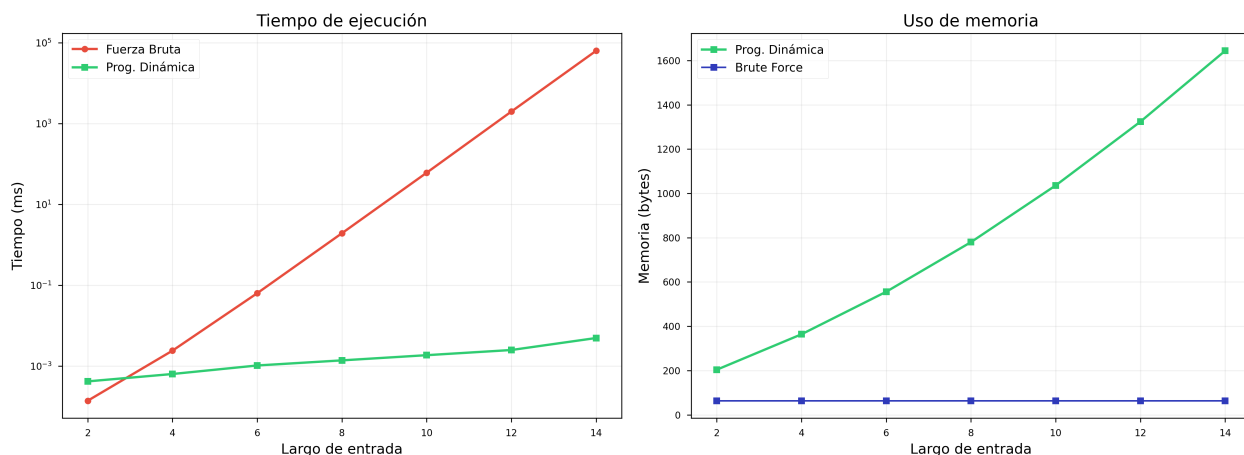


Figura 2: Cadenas de pares transpuestos

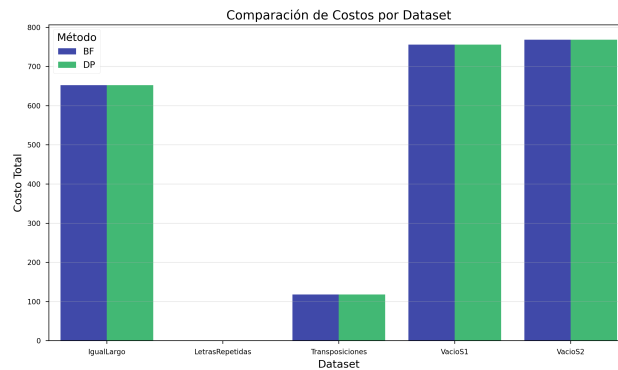


Figura 3: Costos de edicion

Analizando las figuras 2 y 3 se tiene que las transposiciones juegan un papel relevante puesto que la inclusion de estas aumentan la complejidad, pero esto es por una buena razon, ya que estas bajan el costo de edicion de las cadenas.

Análisis de Rendimiento - Dataset: LetrasRepetidas

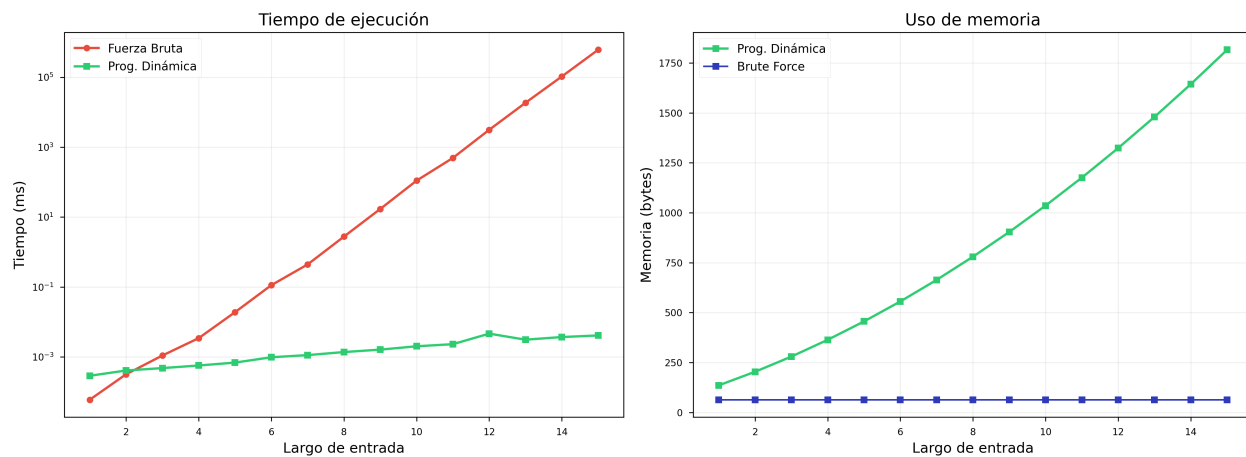


Figura 4: Cadenas de letras repetidas de largo creciente

Las cadenas repetidas es importante estudiarlas, puesto que estas representan el peor caso en temas de complejidad para fuerza bruta, ya que las 4 operaciones se realizaran siempre.

Análisis de Casos Cadenas Vacías

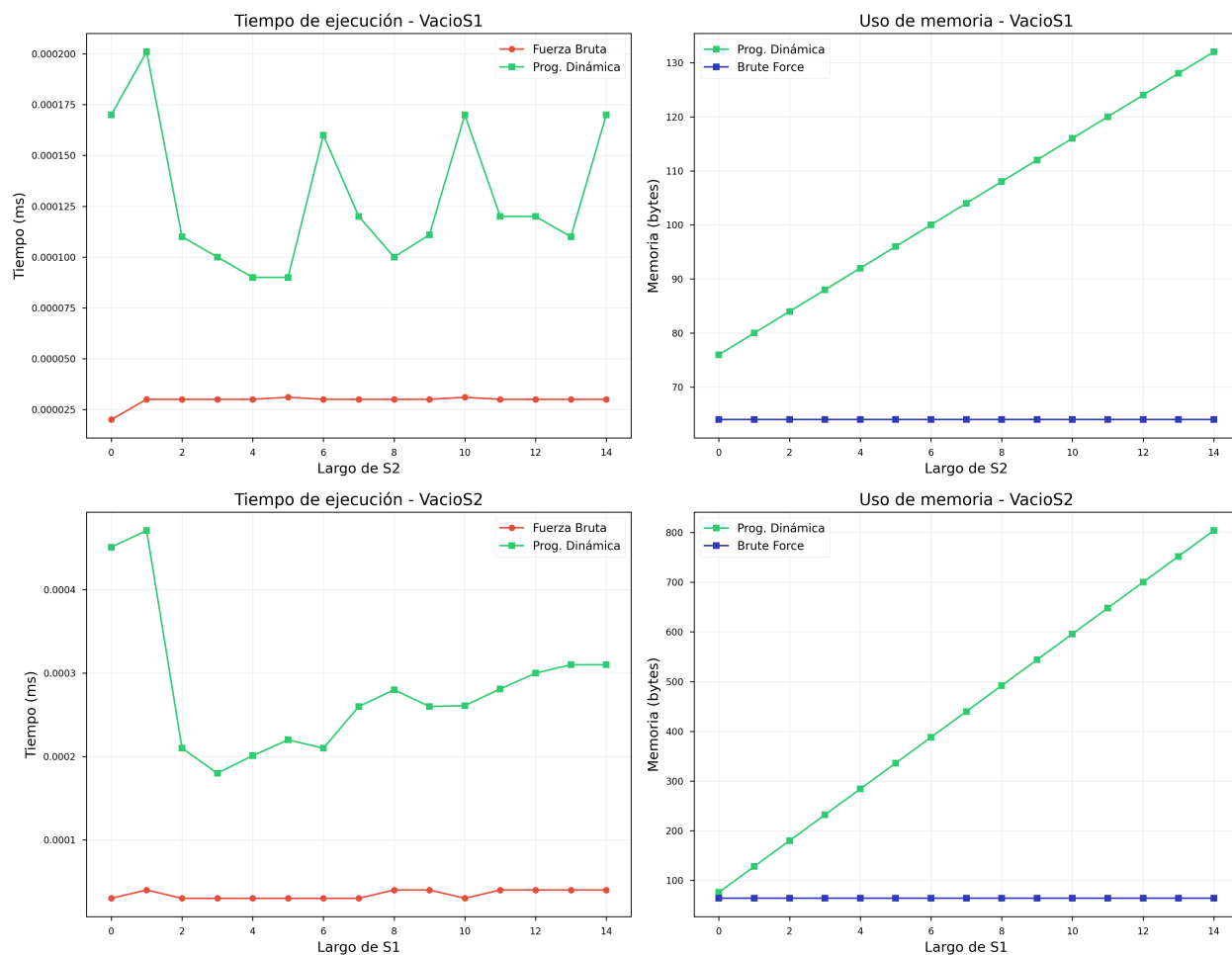


Figura 5: Casos Cadenas Vacías

Las cadenas vacías representan el mejor caso para fuerza bruta, ya que realiza trabajo lineal en vez de exponencial, por esta misma razón le gana en tiempo de ejecución a programación dinámica.

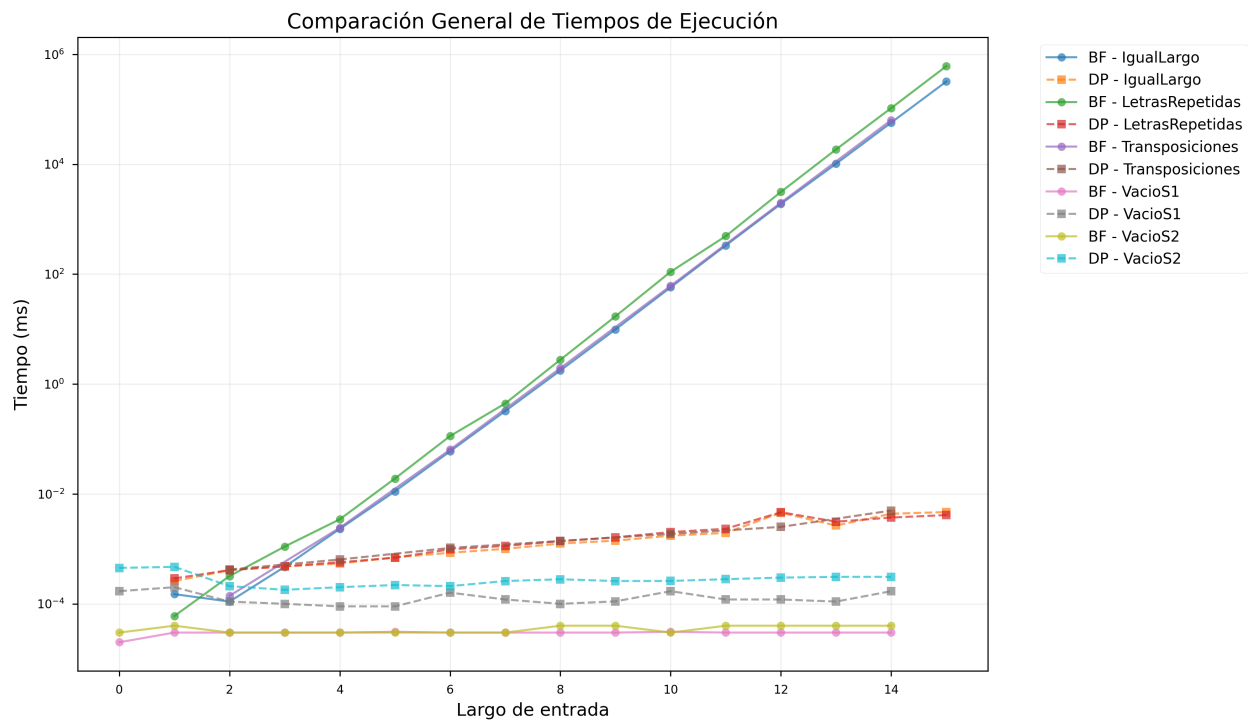


Figura 6: Resumen general de resultados

Finalmente si queremos resolver el problema de la menor distancia, es importante tener en cuenta que el enfoque de fuerza bruta para cadenas vacías puede llegar a ser útil, sin embargo el enfoque de programación dinámica es una buena opción general.

5. Conclusiones

El análisis realizado demuestra que la programación dinámica es una estrategia superior en términos de eficiencia temporal para el cálculo de la distancia de Damerau-Levenshtein, particularmente en casos donde las cadenas tienen longitudes significativas o presentan múltiples posibles operaciones. Al almacenar los resultados intermedios en una tabla de caché, este enfoque reduce drásticamente el número de cálculos redundantes, logrando una complejidad temporal y espacial controlada de $O(n \times m)$.

Por otro lado, el algoritmo de fuerza bruta, aunque limitado por su alta complejidad exponencial en la mayoría de los escenarios, presenta un desempeño competitivo en casos específicos, como cadenas vacías. En estos contextos, su simplicidad lo convierte en una alternativa válida para problemas de menor escala, donde la carga de memoria es un factor crítico.

Este estudio destaca la importancia de seleccionar el enfoque algorítmico adecuado según las características del problema. Además, subraya cómo el análisis detallado de los costos individuales de las operaciones puede influir en la optimización y aplicabilidad de los algoritmos. En general, este trabajo contribuye al entendimiento de la interacción entre complejidad computacional y recursos en el diseño de algoritmos para problemas de edición de cadenas.

6. Condiciones de entrega

- La tarea se realizará **individualmente** (esto es grupos de una persona), sin excepciones.
- La entrega debe realizarse vía <http://aula.usm.cl> en un **tarball** en el área designada al efecto, en el formato `tarea-2 y 3-rol.tar.gz` (rol con dígito verificador y sin guión).
Dicho **tarball** debe contener las fuentes en \LaTeX (al menos `tarea-2 y 3.tex`) de la parte escrita de su entrega, además de un archivo `tarea-2 y 3.pdf`, correspondiente a la compilación de esas fuentes.
- Si se utiliza algún código, idea, o contenido extraído de otra fuente, este **debe** ser citado en el lugar exacto donde se utilice, en lugar de mencionarlo al final del informe.
- Asegúrese que todas sus entregas tengan sus datos completos: número de la tarea, ramo, semestre, nombre y rol. Puede incluirlas como comentarios en sus fuentes \LaTeX (en \TeX comentarios son desde % hasta el final de la línea) o en posibles programas. Anótese como autor de los textos.
- Si usa material adicional al discutido en clases, detállelo. Agregue información suficiente para ubicar ese material (en caso de no tratarse de discusiones con compañeros de curso u otras personas).
- No modifique `preamble.tex`, `tarea_main.tex`, `condiciones.tex`, estructura de directorios, nombres de archivos, configuración del documento, etc. Sólo agregue texto, imágenes, tablas, código, etc. En el código fuente de su informe, no agregue paquetes, ni archivos `.tex` (a excepción de que agregue archivos en `/tikz`, donde puede agregar archivos `.tex` con las fuentes de gráficos en `TikZ`).
- La fecha límite de entrega es el día **10 de noviembre de 2024**.

NO SE ACEPTARÁN TAREAS FUERA DE PLAZO.

- Nos reservamos el derecho de llamar a interrogación sobre algunas de las tareas entregadas. En tal caso, la nota de la tarea será la obtenida en la interrogación.

NO PRESENTARSE A UN LLAMADO A INTERROGACIÓN SIN JUSTIFICACIÓN PREVIA SIGNIFICA AUTOMÁTICAMENTE NOTA 0.

A. Apéndice 1

[Github](#)