

# Proyecto - ADA

Andrés Lostaunau

Junio 2021

## 1 Problema del Trie eficiente

### 1.1 Introducción

El Trie nos presenta una estructura de datos parecida a la de un árbol parejo, cada arista tiene un valor de un carácter que pertenece a un alfabeto  $\Sigma$  y además cumple con una regla específica (en este caso que todos los aristas de un mismo nodo padre están en orden alfabético). Un ptrie cumple con los mismos requerimientos solo que todas las hojas de la misma altura pueden cambiar de nivel en cualquier nivel al igual que una permutación  $p=h!$ . Y un S-ptrie además de lo indicado recibe un conjunto de cadenas  $S$  para construir la estructura. El problema requerido es ver cual de las permutaciones posibles del S-ptrie tiene la menor cantidad de aristas.

Para resolver el problema primero es necesario identificar el causante. En este caso el enunciado es muy claro, son las posibles permutaciones. Dejando de lado este problema debemos tener bien definido el resto de fragmentos del algoritmo implementado.

### 1.2 Implementación General

Lo primero a definir es que el código ha sido implementado en C++. La estructura tomada del S-ptrie se interpretó como una matriz de caracteres del tamaño  $m * n$ , donde  $n$  es igual número de cadenas de  $S$  y  $m$  es igual a la longitud de cada cadena. Cada uno de los caracteres representa a un arista del S-ptrie siendo que todos los caracteres en el mismo eje horizontal se encuentran en la misma altura. Los caracteres @ representan a la ausencia de arista. Si un arista es precedido por un arista ausente, entonces el arista pertenece al predecesor del arista anterior. El número de aristas es igual a  $n * m - r$ , donde  $r$  es igual al número de aristas ausentes (Figure 1).

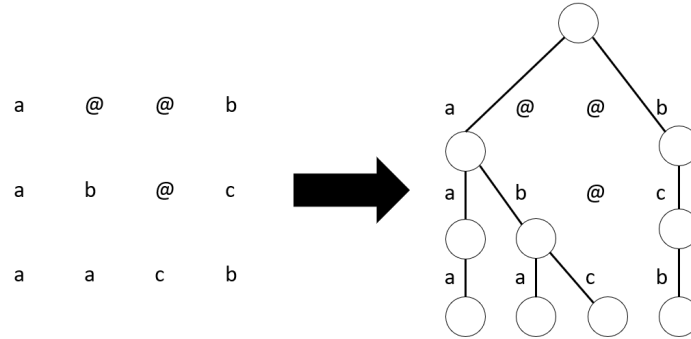


Figure 1: Representación del S-trie

### 1.2.1 Transpose Algorithm

El input original es recibido como un vector de strings con las dimensiones  $n*m$ , inverso a las dimensiones necesitadas para construir un S-trie, para esto se implementa la función *transpose(input)* que sirve para obtener la matriz transpuesta del input.

---

```

1 vector<string> transpose(vector<string> input){
2     vector<string> t_matrix;
3     for(int i = 0; i < input[0].size(); i++){
4         t_matrix.emplace_back("");
5         for(int j = 0; j < input.size(); j++){
6             t_matrix[i].push_back(input[j][i]);
7         }
8     }
9     return t_matrix;
10 }
```

---

El valor de  $m$  está representado como *input.size()*, al igual que  $n$  está representado por *input[0].size()*, esto se da múltiples veces en el código. El algoritmo son dos iteraciones anidadas, por ende el tiempo de ejecución del algoritmo es de  $O(nm)$ .

### 1.2.2 Build Algorithm

---

```

1 int build_strie(vector<string>* input){
2     int removed_edge_counter = 0;
3     char current_char = ' ';
4     for(int i = 0; i < input->size(); i++){
5         for(int j = 0; j < input[0].size(); j++){
6             if((*input)[i][j] != current_char){
7                 current_char = (*input)[i][j];
```

```

8         }else if(i - 1 < 0) {
9             (*input)[i][j] = '@';
10            removed_edge_counter++;
11        }else if((*input)[i - 1][j] == '@') {
12            (*input)[i][j] = '@';
13            removed_edge_counter++;
14        }
15    }
16    current_char=' ';
17 }
18 return removed_edge_counter;
19 }

```

---

El algoritmo de construcción del S-ptrie es algo interesante porque se optó por una estructura de datos más liviana a la propuesta por el enunciado, pero igual de funcional. Al usar una matriz como estructura de datos el espacio que toma se reduce de  $O(mn\Sigma)$  a  $O(mn)$  y el algoritmo de construcción como se ve aquí también comparte el mismo tiempo de ejecución. Ya se explicó el como es que se colocan los aristas ausentes. El algoritmo además devuelve el número de aristas ausentes para calcular posteriormente el número de aristas sin necesidad de volver a iterar.

### 1.2.3 Permutación Óptima

Por último es necesario saber cual de las tantas permutaciones posibles es la óptima para resolver el problema. El principal asunto es que este paso no es sencillo o al menos es muy costoso. Para evaluar todos los casos en condiciones normales tendría un costo inalcanzable de  $O(n!)$ , por esta razón más adelante veremos diferentes algoritmos que buscan poder realizar esta tarea con una facilidad mucho mayor.

## 2 Heurística Voraz

Crear un algoritmo voraz usualmente es una tarea sencilla y encima los algoritmos tienden a ser poco costosos. Observando el problema detenidamente puedes notar un pequeño patrón con respecto al número de aristas y el número de caracteres diferentes en los primeros niveles. Mientras menos caracteres diferentes se encuentren dentro de los primeros niveles el número de aristas será menor. Aquí una pequeña demostración con un  $S=aaa,aab,abc,adb$  donde se tomarán dos instancias donde se puede apreciar este efecto más. (Figure 2)

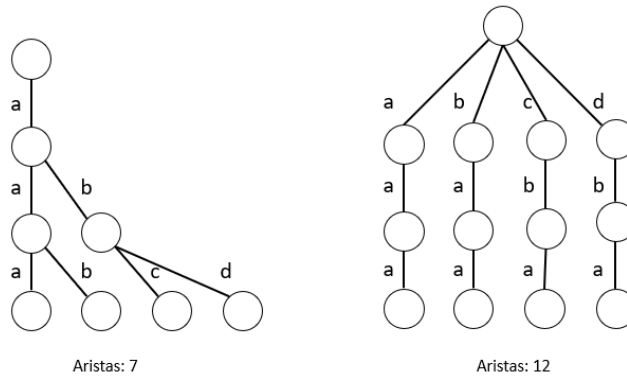


Figure 2: Ejemplo del impacto del orden de los niveles

Bajo este razonamiento es posible crear un algoritmo que ordene la matriz priorizando los niveles con menor repetición primero. Para este algoritmo se han empleado dos partes, la creación de un vector de *score* y realizar un sort.

## 2.1 Implementación

### 2.1.1 Score Vector

---

```

1 vector<int> get_score(vector<string> t_matrix){
2     int score;
3     char actual_char;
4     vector<int> score_vec;
5     for(int i = 0; i < t_matrix.size(); i++){
6         actual_char = ' ';
7         score = 0;
8         for(int j = 0; j < t_matrix[0].size(); j++){
9             if(t_matrix[i][j] != actual_char){
10                 score++;
11                 actual_char=t_matrix[i][j];
12             }
13         }
14         score_vec.emplace_back(score);
15     }
16     return score_vec;
17 }
```

---

La finalidad del vector de score es usarlo como apoyo para el algoritmo de sort. El tiempo de ejecución es de  $O(mn)$

### 2.1.2 Sorting Algorithm

El algoritmo de sort a elección fue un MergeSort personalizado que recibía tanto el vector de score y la matriz, de modo que el algoritmo ordenaba el vector de score y la matriz imitaba los cambios que se realizaban en el vector. Para no llenar de mucho código solo compartiré el algoritmo de Merge.

---

```
1 void merge(vector<int>* score_vec, vector<string>* input, int l, int
    mid, int r){
2     vector<int> lArray, rArray, orderedArray;
3     vector<string> lMatrixArray, rMatrixArray, orderedMatrixArray;
4     for(int i = l; i < mid; i++){
5         lArray.push_back((*score_vec)[i]);
6         lMatrixArray.push_back((*input)[i]);
7     }
8     for(int i = mid; i < r; i++){
9         rArray.push_back((*score_vec)[i]);
10        rMatrixArray.push_back((*input)[i]);
11    }
12    int rIndex = 0, lIndex = 0;
13    while(lIndex != lArray.size() || rIndex != rArray.size()){
14        if(lIndex < lArray.size() && (rIndex == rArray.size() ||
            (lArray[lIndex] < rArray[rIndex]))){
15            orderedArray.emplace_back(lArray[lIndex]);
16            orderedMatrixArray.emplace_back(lMatrixArray[lIndex]);
17            lIndex++;
18        }else{
19            orderedArray.emplace_back(rArray[rIndex]);
20            orderedMatrixArray.emplace_back(rMatrixArray[rIndex]);
21            rIndex++;
22        }
23    }
24    for(int i = l; i < r; i++){
25        (*score_vec)[i] = orderedArray[i-l];
26        (*input)[i] = orderedMatrixArray[i-l];
27    }
28 }
```

---

Ya que es un MergeSort el tiempo de ejecución del algoritmo es de  $O(m \lg(m))$ .

## 2.2 Ejecución

---

```
1 struct Answer{
2     vector<string> sptrie;
3     int edges;
4 };
5
6 Answer greedy_algorithm(vector<string> const input){
7     vector<string> t_matrix = transpose(input);
```

```

8      sort(&t_matrix);
9      int removed_edges = build_strie(&t_matrix);
10
11      Answer ans;
12      ans.edges = (t_matrix.size()*t_matrix[0].size()) - removed_edges;
13      ans.sptrie = t_matrix;
14      return ans;
15  }

```

---

Esta es la ejecución del algoritmo voraz que devuelve tanto el S-ptrie como el número de aristas. El tiempo de ejecución del algoritmo es igual a la suma de los tres fragmentos mencionados previamente de  $O(mn)$ ,  $O(mlg(m))$  y  $O(mn)$  dejando al algoritmo con un tiempo de ejecución de  $O(mn + mlg(m))$ . Se realizó una prueba con un  $S=aaa,baa,bac,cbb$ , siendo el resultado el de Figure 3.

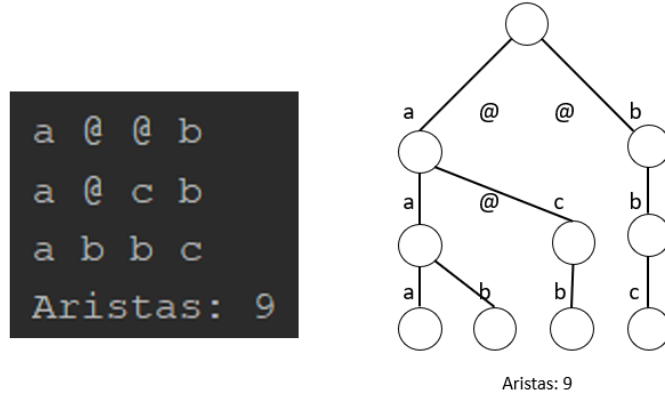


Figure 3: Ejemplo del impacto del orden de los niveles

El resultado es el mínimo entre las permutaciones posibles, por ende es correcto para este caso.

## 2.3 Demostración

El algoritmo que determinaba la permutación adecuada fue presentado como un algoritmo voraz correcto, sin embargo para poder determinar su exactitud es necesario demostrar la propiedad de la elección voraz y la propiedad de la subestructura óptima.

Antes de entrar a las demostraciones específicas es necesario explicar la ley del score. Si tenemos un vector de score  $V$  con  $n$  niveles y un vector de aristas  $A$  con  $n$  niveles que representan el número de aristas entre el nivel  $n-1$  y  $n$  en el S-ptrie podemos decir que siempre se cumplirá que para todo  $v_i \in V$  tiene un  $a_i \in A$  donde  $v_i \leq a_i$ ,  $a_{i-1} \leq a_i$  y  $a_i \leq n$  a excepción de  $a_1 = v_1$ .

### 2.3.1 Elección Voraz

**Lema:** Si existe un vector de score ordenado  $V$  donde  $V=\{v_1, v_2, \dots, v_n\}$  y  $v_i \in V, v_{i-1} < v_i$ , debe ser el óptimo.

**Prueba:** Si existe un vector  $X$  donde con  $n$  valores diferentes y con el orden  $x_i \in X, 1 \leq x_i \leq n, x_{i-1} < x_i$  será óptimo.

Sabiendo que  $a_1 = x_1 = 1$  y que  $a_{i-1} \leq a_i$  es posible determinar que  $x_1 \leq a_2$  y considerando que  $v_i \leq a_i$  también que  $x_2 \leq a_2$  pero como sabemos que  $x_1 = 1$  entonces obtenemos que  $x_1 < x_2$  y lo mismo se puede aplicar hasta  $x_n$  llegando al final a la forma óptima donde  $x_1 < x_2 < \dots < x_n$

### 2.3.2 Subestructura Óptima

**Lema:** Un vector  $V$  que contiene dos o más elementos  $v_i = v_j$  no cumple con la elección voraz.

**Proof:** Tenemos un vector  $X$  que contiene a un par de elementos  $x_i = x_j$ . Si asumimos que  $j = i+1$  tendremos que  $x_i = x_i + 1, x_i \leq a_i, x_{i+1} \leq a_{i+1}$  y  $a_i \leq a_{i+1}$  lo que nos dejaría en  $x_i \leq a_i \leq a_{i+1}$ . Como no es posible identificar si  $a_i < a_{i+1}$  o  $a_i = a_{i+1}$  y eso influye en el valor óptimo, entonces  $X$  no es pertenece al conjunto de respuestas cuando existe un  $x_i = x_j$ .

## 3 Recurrencia

El problema del algoritmo presentado anteriormente es que al momento del ordenamiento este no distingue de niveles con el mismo score por ende si el orden de niveles con el mismo score fuese relevante el algoritmo no aseguraría su respuesta. Es muy poco probable que todos los niveles tengan el mismo puntaje así que no se apostar como un algoritmo confiable para resolver este problema. Observando nuevamente el problema es posible darle un nuevo enfoque, un enfoque recurrente.

Analizando nuevamente la matriz transpuesta es posible notar otro comportamiento, cada vez que una cadena no se relaciona con otra cadena en un nivel se da que por el resto de niveles la cadena no va a relacionarse, por ende siempre va a generar una arista. Entonces si es que es posible separar a todas las cadenas que ni bien deje de relacionarse. Mientras que una cadena se siga relacionando con otra sigue siendo tomada como candidata y se sigue tomando en consideración. (Figure 4)

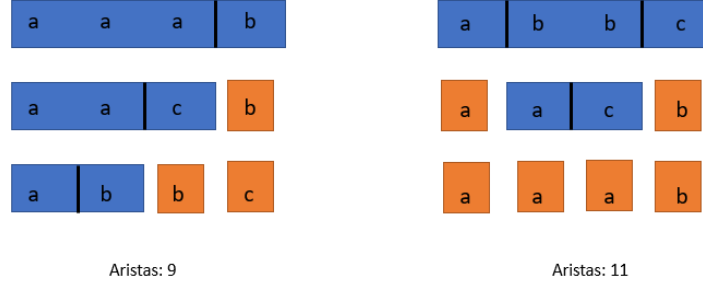


Figure 4: Ejemplo del conteo de aristas

Entonces con esta interpretación podemos definir un subproblema de la siguiente forma:

$$OPT(m, i, j) = \begin{cases} 0 & , m = 0 \\ m & , i = j \\ (OPT(m-1, i, a) + 1) + \\ (OPT(m-1, a+1, b) + 1) \\ + \dots + \\ (OPT(m-1, z+1, j) + 1) & , cc \end{cases} \quad (1)$$

Para interpretarlo con más facilidad la variable  $m$  representa la altura dentro del S-ptrie, la variable  $i$  y  $t$  son el índice inferior y superior respectivamente y las variables  $a, b, \dots, z$  representan los índices en donde ocurre una división. Con esta estructura podemos crear un grafo acíclico dirigido que permita una mejor visualización de lo que pasa.

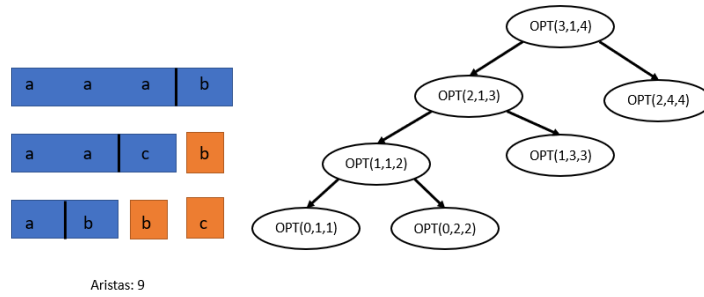


Figure 5: Ejemplo con  $S=\{aaa,aab,acb,bbc\}$

Aquí se puede apreciar como el problema se descompone subproblema por subproblema de modo que devuelve el resultado esperado.