



UNIVERSIDAD
DE LA REPÚBLICA
URUGUAY



FACULTAD DE
INGENIERÍA
UDELAR

Algoritmos eficientes para la construcción de conjuntos óptimos de contextos bi–direccionales

Fernando Arsenio Fernández Barreiro

Tesis de Maestría presentada a la Facultad de Ingeniería de la Universidad de la República
en cumplimiento parcial de los requerimientos para la obtención del título de Magister en
Informática.

Tutor

Marcelo Weinberger

Tribunal

Gonzalo Navarro
Álvaro Martín
Pedro Piñeyro

Montevideo, Uruguay
octubre de 2018

Tesis de maestría
Algoritmos eficientes para la construcción de
conjuntos óptimos de contextos
bi–direccionales

Fernando Fernández Barreiro
fefernan@fing.edu.uy

Director de Tesis: Dr. Marcelo Weinberger
Center for Science of Information
marcwein@ieee.org

Director Académico: Dr. Alberto Pardo
Universidad de la República, Montevideo, Uruguay
pardo@fing.edu.uy

Octubre, 2018

Resumen

Resumen La técnica de modelados por contextos usada en problemas de decisión secuencial, como compresión de datos y otros, descompone una secuencia dada entre un conjunto de subsecuencias que son tratadas independientemente. La descomposición se basa en la ocurrencia de ciertas cadenas de símbolos que son elementos de un conjunto finito de contextos. Cada contexto tiene asociado un peso dado por una pérdida numérica generada al procesar la subsecuencia. Por tal motivo es necesario encontrar un conjunto de contextos óptimo (que minimice la pérdida total). En esta tesis los contextos están formados por pares de cadenas y se denominan contextos bi-direccionales.

La contribución de este trabajo es presentar estructuras de datos y algoritmos computacionalmente eficientes para encontrar un conjunto óptimo de contextos cuando la longitud de los contextos no está acotada.

Se demuestra que el conjunto de todos estos contextos bi-direccionales se puede partir en clases de equivalencia y alcanza con procesar un elemento representativo de cada clase.

La estructura presentada es una generalización de árboles de sufijos compactos. Es un grafo enraizado, dirigido, acíclico, donde cada vértice está etiquetado por el elemento representativo de una clase. Recíprocamente a cada clase le corresponde un vértice de este grafo. Las aristas dirigidas (al igual que en los árboles compactos de sufijos) están etiquetados por cadenas de símbolos, con las que se extiende el contexto representado por el vértice origen. Algunos vértices no tienen aristas salientes. Cada uno de los otros vértices tiene dos conjuntos de hijos (cada uno de los cuales extiende una de las direcciones de los contextos).

También se describe un algoritmo que reduce la complejidad de orden cúbico a orden cuadrático (en el largo de la cadena). Para ello se hace uso de árboles de sufijos y el cálculo del ancestro común más profundo de varios nodos de un árbol.

Palabras clave: Modelado por contexto, contextos bi-direccionales, árboles de contexto, compresión, predicción, eliminación de ruido, árbol de sufijos, ancestro común más profundo, grafos.

Índice general

1. Introducción	5
1.1. Motivación del problema	6
1.2. Definición del problema	8
2. Definiciones fundamentales	11
3. Fundamentos sobre conjuntos óptimos de contextos	18
3.1. Contextos uni–direccionales	18
3.2. Contextos bi–direccionales	22
4. Contextos Representativos	28
4.1. Caso uni–direccional	29
4.2. Generalización Bi–direccional	32
4.3. Contextos Representativos	33
4.3.1. Conjunto óptimo	38
5. Subárboles Compactos	41
5.1. Árboles	41
5.1.1. Subárboles Compactos	42
5.1.2. LCA	43
5.1.3. LCA generalizado	44
5.1.4. Árbol soporte	47
5.1.5. Unión de subárboles	47
5.2. Búsqueda de contextos representativos	51
5.2.1. Obtención de los contextos representativos	51
5.2.2. Cantidad de contextos representativos	54
6. Grafos Compactos de Contextos Bi–direccionales	57
6.1. Definición	58
6.2. Construcción del CBDCG	61
6.3. Podado	66

6.4. Complejidad del Algoritmo	67
6.5. Datos experimentales	68
6.5.1. Consideraciones prácticas	68
6.5.2. Resultados experimentales	70
7. Conclusión	75
Bibliografía	78
A. Cálculo del $1ca$	81

Capítulo 1

Introducción

En esta tesis se presenta una estructura de datos para manejar contextos bidireccionales que es una generalización de los árboles de sufijos, y un algoritmo eficiente para construirla. Es una generalización de árboles, en donde cada nodo representa un contexto y tiene dos conjuntos de hijos que representan extensiones de sus contextos. Resulta ser un caso particular de grafo enraizado, dirigido y acíclico.

La motivación es contribuir en la técnica de *modelado por contextos* usada en problemas como *compresión, predicción* y eliminación de ruido (*denoising*). En estas aplicaciones se procesan secuencias (también llamadas cadenas) de símbolos y para cada posición de esas secuencias se tiene en cuenta sus *contextos*, que son subsecuencias que rodean a la posición. Una subsecuencia de una cadena es una secuencia formada por los símbolos que están en un conjunto de posiciones de la secuencia original. Por ejemplo, dada la secuencia *abcbab*, la subsecuencia correspondiente a las posiciones $\{1, 4\}$ es *aa*. Una subcadena es una subsecuencia determinada por posiciones consecutivas de la secuencia original, por ejemplo *bca*.

La secuencia a procesar se descompone en un conjunto de subsecuencias que son tratadas independientemente. La descomposición se basa en la ocurrencia de contextos en las posiciones de la secuencia. Los contextos pertenecen a un conjunto finito que debe cumplir las propiedades de ser disjunto y exhaustivo, lo cual implica que a cada posición de la secuencia le corresponde exactamente un contexto del conjunto. El procesamiento de la secuencia asigna a cada una de las subsecuencias un costo (pérdida numérica) que está determinado por la cantidad de ocurrencias de cada símbolo. Por extensión, al contexto correspondiente a la subsecuencia se le puede asociar ese costo. Esto motiva a encontrar un conjunto de contextos que sea óptimo (esto es, que minimice el costo total) entre aquellos que cumplan las propiedades requeridas.

Es posible que el contexto esté formado por varios componentes. Otro aspecto

que es de interés es la longitud de los contextos. Puede estar acotada o no y puede ser la misma en todos los contextos o no. Si los contextos tienen varios componentes se puede exigir que el largo de los componentes sea el mismo. En el caso tratado en este documento los contextos tienen dos componentes, las longitudes pueden ser diferentes entre contextos y entre componentes de un mismo contexto y no se acota la longitud que pueden tomar. Esto último genera un problema de complejidad computacional y es lo que aquí se resuelve. Para lograrlo lo que se va a hacer es particionar el conjunto de contextos posibles en clases, de tal manera que los elementos de la misma clase tengan propiedades idénticas. Y se va a procesar un único contexto de cada clase.

1.1. Motivación del problema

En el modelado por contextos [1, 2, 3] se procesa una secuencia descomponiéndola en subsecuencias y procesando cada subsecuencia de manera independiente. El procesamiento de cada subsecuencia genera un costo y se busca una descomposición que minimice el costo total, que es la suma de los costos por cada subsecuencia. Las subsecuencias se forman seleccionando posiciones que tienen un contexto en común. Como ejemplo supongamos que los contextos de una posición son las subsecuencias, leídas de derecha a izquierda, que la preceden. De manera provisoria, excluimos la posición 1, ya que no hay subsecuencias que la preceden. Para resolver esto, en la Subsección 3.1 se agregará un símbolo marcador en la posición 0, que no ocurre en el resto de la secuencia. Sea $abcdbccz$ la secuencia procesada. Algunos contextos posibles para la cuarta posición son cb y cba . El primero de estos contextos también precede a la séptima posición. Por lo tanto la subsecuencia asociada a cb es la formada por los elementos que están en las posiciones 4 y 7, dc . Por otro lado, el contexto cba sólo precede la cuarta posición, por lo que le corresponde la subsecuencia d . Entonces, lo que hay que decidir es cómo se elige el conjunto de contextos que determinan las subsecuencias.

El conjunto debe satisfacer dos propiedades: *exhaustiva* y *disjunta*, (que se verán en el Capítulo 3), que son condición para que cada posición de la secuencia procesada corresponda a exactamente una de las subsecuencias. A un conjunto que cumpla esas propiedades se le dice *conjunto válido*.

Volviendo al ejemplo anterior, la secuencia $abcdbccz$, supongamos que cba es uno de los contextos. Precede a la cuarta posición pero no a la séptima. Para que sea un elemento de un conjunto válido tiene que haber un contexto que preceda a la séptima posición, y no puede ser c o cb porque estos preceden también a la cuarta posición. Un contexto posible es cbd , ya que no precede a ninguna otra posición de la secuencia.

Un par de ejemplos de conjuntos válidos para la secuencia $abcdbccz$ son

$\{a, b, cba, cbd, cc, d\}$ y $\{a, ba, bd, cb, cc, d\}$.

Estos son conjuntos válidos específicos de la secuencia tratada, relativos a ella. Se debe tener en cuenta que en la bibliografía lo habitual es definir los conjuntos válidos para *cualquier* secuencia. En el ejemplo, si el alfabeto fuese $\{a, b, c, d, z\}$ un posible conjunto válido de contextos sería $\{a, b, ca, cba, cbb, cbc, cbd, cbz, cc, cd, cz, d, z\}$. Con esta definición de validez cada posición de *cualquier* secuencia es precedida por algún contexto del conjunto (recordamos que la secuencia se lee de derecha a izquierda), aunque puede haber contextos que no precedan ninguna posición.

Entre los conjuntos válidos, aquellos que generan un conjunto de subsecuencias para los cuales el costo total es mínimo se llaman *conjuntos óptimos*. Como ejemplo ilustrativo, sin relación con una aplicación real, supongamos que a las secuencias de longitud uno se asocia un costo de valor 2, mientras que a las secuencias de dos símbolos corresponde un costo de valor 3 si ambos símbolos son iguales y de valor 5 si los símbolos son diferentes. Como se vio, el contexto cb precede a las posiciones cuarta y séptima y determina la subsecuencia dc . El costo que corresponde es 5. Pero para cubrir las mismas posiciones también se pueden elegir los contextos cba y cbd , los cuales determinan subsecuencias de longitud 1, por lo que el costo total asociado a esas posiciones sería 4. Entonces, un conjunto válido que incluya el contexto cb no puede ser óptimo, ya que “partiendo” ese contexto en el par $\{cba, cbd\}$ se obtiene un conjunto con el que se asocia un costo menor. Lo opuesto sucede al considerar las posiciones tercera y sexta. Se puede elegir los contextos ba y bd que determinan subsecuencias de longitud 1 y, por lo tanto, con un costo asociado de valor 4. Pero también se puede elegir el contexto b . La subsecuencia determinada es cc , y el costo es 3. De donde, si el conjunto válido contiene el par $\{ba, bd\}$, es conveniente “unificar” esos contextos en el contexto b . El costo de $\{a, b, cba, cbd, cc, d\}$ es 13 y el de $\{a, ba, bd, cb, cc, d\}$ es 15.

Es conocido que un conjunto válido se puede representar con las hojas de un árbol m -ario siendo m la cantidad de símbolos del alfabeto sobre el que se define la secuencia procesada [2]. Si la longitud de los contextos es acotada por k , se conocen algoritmos que, podando el árbol m -ario perfecto de altura k , obtienen un árbol que representa un conjunto óptimo. A estos árboles se les denomina *árboles óptimos*. Si no se acota la longitud de los contextos, se conoce un algoritmo que también permite encontrar eficientemente un conjunto óptimo [4].

Se ha formalizado la definición de conjunto válido cuando los contextos están formados por pares de secuencias, no necesariamente de la misma longitud, dando la generalización de las propiedades exhaustiva y disjunta [5, 6, 7]. Se demostró que un conjunto válido en el llamado caso *bi-direccional* también puede representarse con las hojas de un árbol m -ario. Y se describió un algoritmo que encuentra un conjunto óptimo si se mantiene acotada la longitud de los contextos. Este es

un algoritmo que empieza en los contextos más largos y mediante programación dinámica va podando una estructura que está implícita y es un grafo.

La demostración de que en el caso bi–direccional un conjunto válido se representa mediante hojas de un árbol no es una simple formalización que generaliza el caso uni–direccional, ya que se ha demostrado que esto mismo no se cumple cuando la cantidad de componentes de los contextos es mayor a 2 [6, 7].

1.2. Definición del problema

En principio se puede suponer que establecer un límite en la longitud de los contextos reduce la información que se obtiene de la secuencia. Lo que acá nos proponemos es eliminar la cota en la longitud de los contextos. En realidad hay una cota implícita que es la longitud de la secuencia a procesar. Pero aplicar el algoritmo mencionado con esa cota implica complejidad de orden cúbico en la longitud de la secuencia. Entonces, el propósito de este trabajo es diseñar un algoritmo que encuentre un conjunto óptimo de contextos bi–direccionales cuando no se acota la longitud de los contextos, y que en el peor caso tenga complejidad polinómica de orden cuadrático.

La solución pasa por observar que para algunos contextos la información relevante al problema es la misma que para otros. Entonces, para cada conjunto de contextos con idéntica información, lo que se debe hacer es encontrar sólo uno de ellos, que los represente a todos. Lo que sucede con cada uno de los contextos de ese conjunto es que en cada posición de la secuencia en que ocurre, también ocurre su *contexto representativo*, del cual es *prefijo*. Esto se desarrolla en el Capítulo 4.

Se propondrá un algoritmo que encuentre estos contextos representativos. Además, se demostrará que cada contexto que ocurre en la secuencia de entrada está representado por uno de los contextos encontrados por el algoritmo. Parte del trabajo consiste en volver a definir árboles y árboles óptimos en los que los nodos no sean todos los contextos posibles, ni siquiera todos los que ocurren, sino sólo aquellos que representan a otros. También es necesario definir nuevas estructuras, grafos compactos con raíz, dirigidos y acíclicos, cuyos nodos corresponden a esos contextos representativos.

Tal como en la versión uni–direccional [7], el algoritmo empieza con los contextos de longitud máxima. Los contextos están formados por un par de componentes, que llamamos izquierdo y derecho. Todos los contextos izquierdos (derechos) son subcadenas de una cadena denominada *cadena de formación de contextos izquierdos (derechos)*. A cada posición de la secuencia procesada corresponde un sufijo de esta cadena de formación de contextos. Cada sufijo es la cadena formadora de contextos de su posición correspondiente. Los contextos son los prefijos de estos sufijos. Los contextos bi–direccionales de longitud máxima

son los que alcanzan el final de la cadena de formación de contextos correspondiente en cada una de los dos componentes. A los contextos de longitud máxima les llamamos *bi-sufijos*, porque sus dos componentes son sufijos de las cadenas de formación de contextos. La información relevante de un contexto, mencionada anteriormente, está determinada por los bi-sufijos de los cuales es prefijo. Hay una correspondencia entre los bi-sufijos y las posiciones de la secuencia de entrada, ya que cada posición define un contexto de longitud máxima.

En el caso uni-direccional un contexto se dice *ramificable* si ocurre en la cadena de formación de contextos al menos dos veces y el símbolo siguiente es diferente en ambas ocurrencias. Al plantear el ejemplo de la secuencia $abcdbccz$ se vio que el contexto cb ocurre dos veces. En ese ejemplo la cadena de formación de contextos es $ccbdcb$. En esta cadena cb es ramificable porque una de las veces que ocurre es prefijo de cba y la otra de cbd . En el árbol de contextos que corresponde a la secuencia, hay un nodo que representa el contexto cb del cual salen dos ramas. En el caso bi-direccional le diremos *doblemente ramificable* o *bi-ramificable* si es ramificable en ambas direcciones. Veremos que los contextos representativos son los bi-sufijos y los bi-ramificables.

Mientras que es fácil encontrar los bi-sufijos, la complejidad del algoritmo aquí desarrollado consiste en encontrar los bi-ramificables. Todos los contextos izquierdos se pueden representar con un *árbol de sufijos*, y los contextos derechos con otro. Los contextos bi-direccionales representativos son un subconjunto del producto cartesiano de los nodos de esos árboles. Entre las hojas de ambos árboles hay una correspondencia uno a uno que determina los bi-sufijos y se obtiene de manera trivial. Empezando con esos contextos de longitud máxima se recorre desde las hojas hacia la raíz uno de los árboles, encontrando para cada nodo del árbol recorrido aquellos nodos del otro árbol con los que forma un contexto representativo. La principal operación que se usa es el cálculo del prefijo común más largo de varias secuencias. El costo de esta operación es proporcional a la longitud de las cadenas. Pero, si las cadenas son nodos de un árbol, como sucede en este caso, esa operación se puede transformar en la obtención del ancestro común más profundo o *lca* (por su denominación en inglés “lowest common ancestor”), que se puede hacer en tiempo constante después de un preprocesamiento del árbol.

Además de encontrar los nodos que corresponden a los contextos representativos hay que estructurarlos según relaciones padre-hijo, porque el objetivo es podar la estructura construida para obtener un árbol óptimo que está contenido en ella.

Las contribuciones fundamentales de esta tesis fueron presentadas en seminarios y conferencias [8, 9].

En el Capítulo 2 se definen conceptos básicos que se usarán en el resto del documento. En el Capítulo 3 se describe el problema cuando se fija una cota a la longitud de los contextos, tanto en el caso uni como en el bi-direccional. En el Capítulo 4 se trata el caso no acotado. Los algoritmos y las estructuras se describen en el Capítulo 6. Previamente, en el Capítulo 5 se describe la versión uni-direccional del algoritmo principal. En el Apéndice A se describe la forma de resolver el cálculo del LCA.

Capítulo 2

Definiciones fundamentales

En este capítulo se presentan la definición y la notación de los principales conceptos con los que se tratará en el resto del documento.

Cadenas

A lo largo de todo el documento usaremos un *alfabeto finito* \mathcal{A} . Con \mathcal{A}^* representamos el conjunto de todas las cadenas finitas de símbolos de \mathcal{A} , incluyendo la cadena de longitud nula, denotada por ϵ . Con \mathcal{A}^+ se simboliza el conjunto de cadenas finitas de longitud no nula. \mathcal{A}^k es el conjunto de cadenas de longitud k y $\mathcal{A}^{\geq k}$ el de cadenas de longitud mayor o igual a k .

En algunas ocasiones se trabaja con pares ordenados de cadenas finitas. En general las dos cadenas del par no tienen por qué tener la misma longitud. Al primer componente le llamamos izquierdo y al segundo derecho. En este caso se dice que hay dos direcciones, izquierda y derecha, y lo definimos como caso *bi-direccional*. Por compatibilidad con esto, al caso simple le decimos también *uni-direccional*. Si el par de cadenas es (ϵ, ϵ) , con un abuso de notación y asumiendo que no hay ambigüedad, lo denotaremos también por ϵ .

Siempre supondremos que procesamos una *secuencia* o *cadena* \mathbf{x} de longitud n , (x_1, x_2, \dots, x_n) . En el caso uni-direccional los elementos de \mathbf{x} pertenecen a \mathcal{A} . En el caso bi-direccional hay aplicaciones, como compresión de señales de audio que provienen de dos canales, en que pertenecen a $\mathcal{A} \times \mathcal{A}$. Pero en otras, por ejemplo la eliminación de ruido, son elementos de \mathcal{A} ; a pesar de eso corresponden al caso bi-direccional ya que para cada posición se tienen en cuenta dos conjuntos de símbolos, los que están antes y los que están después. Generalizando, le llamamos \mathcal{X} al conjunto al que pertenecen los elementos de \mathbf{x} , ya sea \mathcal{A} o $\mathcal{A} \times \mathcal{A}$. La secuencia \mathbf{x} pertenece a \mathcal{X}^* , conjunto de secuencias finitas de elementos de \mathcal{X} . Además definimos \mathcal{X}^* para denotar el conjunto de contextos. En el caso uni-direccional es igual a \mathcal{A}^* y a \mathcal{X}^* . Pero en el caso bi-direccional es $\mathcal{A}^* \times \mathcal{A}^*$, tanto para las aplicaciones en que \mathcal{X} es \mathcal{A} como para aquellas en las que \mathcal{X} es $\mathcal{A} \times \mathcal{A}$.

Cada elemento de \mathbf{x} está en una posición o *índice* perteneciente a $\{1, \dots, n\}$. Dada una secuencia de índices $\mathcal{I} = (i_1, i_2, \dots, i_h)$, con $1 \leq i_1 < i_2 < \dots < i_h \leq n$, la *subsecuencia* de \mathbf{x} indizada por \mathcal{I} es la secuencia formada por los elementos de \mathbf{x} que están en las posiciones de \mathcal{I} , esto es $(x_{i_1}, x_{i_2}, \dots, x_{i_h})$.

Ejemplo 2.1. Supongamos que $\mathcal{A} = \{a, b, c, d, e\}$.

Sea la cadena abcdeabcd. Al conjunto índice $\mathcal{I} = \{2, 4\}$ corresponde la subsecuencia bd.

Para el caso bi–direccional en que los elementos de \mathbf{x} pertenecen a $\mathcal{A} \times \mathcal{A}$, sea $\mathbf{x} = \begin{pmatrix} abcdeabcd \\ dcbaedcba \end{pmatrix}$. Con esta notación en la segunda posición de \mathbf{x} está el símbolo (b, c) de \mathcal{X} . A $\mathcal{I} = \{2, 4\}$ corresponde $\begin{pmatrix} bd \\ ca \end{pmatrix}$.

Una *subcadena* es una subsecuencia definida por un conjunto cuyas posiciones son consecutivas: $(x_i, x_{i+1}, \dots, x_j)$, con $i, j \in \{1, \dots, n\}$. Esta subcadena se denota por \mathbf{x}_i^j . Si $i = j$ la subcadena tiene longitud 1. Si $i > j$ se considera que \mathbf{x}_i^j es ϵ , la cadena vacía.

Para cualquier cadena \mathbf{s} su longitud se denota con $|\mathbf{s}|$. El *reverso* de la cadena \mathbf{s} es la misma cadena pero recorrida de derecha a izquierda y se denota por $\bar{\mathbf{s}}$. Si la secuencia es una subcadena \mathbf{x}_i^j , su reverso es $(x_j, x_{j-1}, \dots, x_i)$, y usaremos la notación ${}^j\mathbf{x}$. En este caso la cadena es vacía si el superíndice es mayor que el subíndice.

Sean $\mathbf{s}, \mathbf{t}, \mathbf{u} \in \mathcal{A}^*$, con $\mathbf{s} = (s_1, s_2, \dots, s_h)$, y $\mathbf{t} = (t_1, t_2, \dots, t_k)$. Se denota $\mathbf{u} = \mathbf{st}$ a la *concatenación* de \mathbf{s} y \mathbf{t} definida por $\mathbf{u} = (s_1, s_2, \dots, s_h, t_1, t_2, \dots, t_k)$. Se dice que \mathbf{s} es *prefijo* de \mathbf{u} , lo que se denota por $\mathbf{s} \preceq \mathbf{u}$. Si $\mathbf{t} \neq \epsilon$ las cadenas \mathbf{s} y \mathbf{u} son diferentes, y la relación se denomina *estricta*, y se denota por $\mathbf{s} \prec \mathbf{u}$. La cadena \mathbf{t} es un *sufijo* de \mathbf{u} , y es sufijo estricto si $\mathbf{s} \neq \epsilon$. La cadena ϵ , es prefijo y sufijo de cualquier cadena. Las relaciones prefijo y sufijo definen un orden parcial entre las cadenas de \mathcal{A}^* .

La notación se extiende a conjuntos. Si \mathcal{C} es un conjunto de cadenas, entonces $\mathbf{s} \preceq \mathcal{C}$ significa que \mathbf{s} es prefijo de cada cadena de \mathcal{C} . Lo mismo vale para la relación estricta.

Las relaciones prefijo y sufijo se extienden a elementos de \mathcal{X}^* . Sean $\mathbf{s} = (s_\ell, s_r)$ y $\mathbf{t} = (t_\ell, t_r)$. Entonces \mathbf{s} es prefijo de \mathbf{t} si y sólo si $s_\ell \preceq t_\ell$ y $s_r \preceq t_r$.

En esta tesis, por comodidad en el lenguaje, tanto para el caso uni–direccional como bi–direccional, si \mathbf{s} es prefijo de \mathbf{t} , decimos también que \mathbf{t} es *extensión* de \mathbf{s} . Entonces en lugar de referirnos al “conjunto de cadenas (o pares de cadenas) de las que \mathbf{s} es prefijo”, nos referimos a “las extensiones de \mathbf{s} ”. Ese conjunto de extensiones lo denotamos con $\mathcal{P}(\mathbf{s})$. O sea

$$\mathcal{P}(\mathbf{s}) = \{\mathbf{t} \in \mathcal{X}^* : \mathbf{s} \preceq \mathbf{t}\}. \quad (2.1)$$

También nos interesará considerar las extensiones de longitud igual a algún parámetro preestablecido. Para el caso bi–direccional, sean k_ℓ, k_r enteros no negativos, y (s_ℓ, s_r) un elemento de \mathcal{X}^* que cumple $|s_\ell| \leq k_\ell$ y $|s_r| \leq k_r$. Denotamos por $\mathcal{P}_{(k_\ell, k_r)}(s_\ell, s_r)$ a las extensiones de (s_ℓ, s_r) tal que su componente izquierdo (resp. derecho) tiene longitud igual a k_ℓ (resp. k_r). Formalmente:

$$\mathcal{P}_{(k_\ell, k_r)}(s_\ell, s_r) = \{(\mathbf{t}_\ell, \mathbf{t}_r) \in \mathcal{A}^{k_\ell} \times \mathcal{A}^{k_r} : (s_\ell, s_r) \preceq (\mathbf{t}_\ell, \mathbf{t}_r)\}. \quad (2.2)$$

El caso uni–direccional es una sencilla simplificación de la Expresión (2.2).

Las definiciones se extienden fácilmente a conjuntos de cadenas: una cadena es extensión de un conjunto de cadenas si es extensión de cada cadena del conjunto.

Definimos *extensiones mínimas* de $s \in \mathcal{A}^*$ a las cadenas sa con $a \in \mathcal{A}$. Para el caso bi–direccional, las *extensiones mínimas por izquierda* (resp. *derecha*) de (s_ℓ, s_r) son las cadenas $(s_\ell a, s_r)$ (resp. $(s_\ell, s_r a)$) con $a \in \mathcal{A}$.

Si dadas dos cadenas, ninguna es prefijo de la otra, sus conjuntos de extensiones son disjuntos. Diremos que las dos cadenas son *disjuntas*.

Contextos

El problema a tratar en esta tesis incluye descomponer una secuencia x en un conjunto finito de subsecuencias, de tal manera que cada posición de x sea índice de exactamente una de las subsecuencias. Es decir, se trata de hacer una partición de $\{1, \dots, n\}$ en clases, cada una de las cuales es el conjunto que indiza una subsecuencia.

Ejemplo 2.2. Sea, otra vez, la cadena $x = abcdeabcd$. Los conjuntos de índices $\{2, 4\}$, $\{1, 6, 7, 8\}$ y $\{3, 5, 9\}$ constituyen una partición de $\{1, \dots, 9\}$. Con esta partición, x es descompuesta en $bd, aabc, ced$.

Una manera de formar las subsecuencias es utilizar contextos. De manera informal, un *contexto* de una posición de x es una secuencia de elementos que están en torno a la posición. La definición se hará rigurosa al precisar el significado de “estar en torno”. Los contextos pueden estar constituidos por una sola cadena, o por dos, y en este caso los llamamos también *contextos bi–direccionales* o *bi–contextos*. En un bi–contexto a una de las cadenas la llamamos *contexto izquierdo* y a la otra *contexto derecho*. A los contextos de una sola cadena les llamamos también *simples* o uni–direccionales.

Dada la cadena x , para cada posición en $\{1, \dots, n\}$ se define su o sus *cadena(s) de formación de contextos*, que son cadenas de símbolos de \mathcal{A} . Los contextos de la posición son subsecuencias de la cadena de formación de contextos. En el caso de bi–contextos habrá dos cadenas de formación de contextos: *cadena de formación de contextos izquierdos*, y *cadena de formación de contextos derechos*. En esta

tesis los contextos de una posición dada serán los prefijos de su correspondiente cadena de formación de contextos. En general, las cadenas de formación de contextos serán subcadenas que preceden o suceden a la posición considerada. O sea, para la posición i pueden ser, si x pertenece a \mathcal{A}^* , las subcadenas ${}_{i-1}^1x$ y x_{i+1}^n . Cuando x pertenece a $(\mathcal{A} \times \mathcal{A})^*$, $x = (x_\ell, x_r)$, pueden ser ${}_{i-1}^1[x_\ell]$ y ${}_{i-1}^1[x_r]$.

Por lo tanto, las cadenas de formación de contextos para cada posición son sufijos de una única cadena, ya sea ${}_{n-1}^1x$ o x_2^n . Entonces, en el caso uni-direccional tenemos una única cadena de formación de contextos que denotamos por y , siendo $y = {}_{n-1}^1x$. En el caso bi-direccional hay una cadena de formación de contextos izquierdos, y_ℓ , y una cadena de formación de contextos derechos, y_r . Cada contexto, ya sea simple, izquierdo o derecho, es una subcadena de su respectiva cadena de formación.

Ejemplo 2.3. Sea x la cadena abcdeabcd. Una posible cadena de formación de contextos de la posición 6 es edcba, es decir la subcadena de x desde la posición 5 ($=6 - 1$) hasta la posición 1. La cadena de formación de contextos de x es $y = cbaedcba$, de la cual edcba es sufijo. Algunos contextos para la posición 6 son ϵ , ed, edcba, todos ellos prefijos de edcba.

Otra posible cadena de formación de contextos de la posición 6 es bcd, la subcadena de x desde 7 ($=6+1$) hasta $n = 9$, que es sufijo de bcdeabcd. En el caso bi-direccional esta cadena es y_r , la cadena de formación de contextos derechos de x , mientras que y_ℓ es cbaedcba. Un posible bi-contexto de la posición 6 es (ed, bcd) .

Aquí se justifica la intuición de hablar de direcciones izquierda y derecha: los contextos izquierdos (resp. derechos) corresponden a cadenas que están a la izquierda (resp. derecha) de la posición.

Ejemplo 2.4. Como ejemplo de secuencia x perteneciente a $\mathcal{A}^* \times \mathcal{A}^*$ tomamos $\begin{pmatrix} abcdeabcd \\ dcbaedcba \end{pmatrix}$. Para la posición 6 los componentes de ${}_{6-1}^1x$ son edcba y eabcd. Estas son las cadenas de formación de contextos para la posición 6. Las cadenas de formación de contextos izquierdo y derecho de x son $y_\ell = cbaedcba$ e $y_r = bcdeabcd$.

Conjunto válido y conjunto óptimo

Sea una cadena x a procesar y s un contexto. Definimos el *conjunto índice* de s y lo denotamos por $\mathcal{I}(s)$ al conjunto de posiciones de x de las cuales s es contexto.

Un *conjunto válido de contextos*, \mathcal{V} , (relativo a una secuencia x dada) es un conjunto de contextos tales que sus conjuntos índices determinan una partición de $\{1, \dots, n\}$. Esto es, $\{\mathcal{I}(s) : s \in \mathcal{V}\}$ es una partición de $\{1, \dots, n\}$.

Es posible que para una secuencia x haya más de un conjunto válido de contextos. Cada contexto determina una subsecuencia, y con cada subsecuencia se asocia un valor numérico que representa un costo, que queda asociado a dicho contexto. Se quiere elegir un conjunto válido de contextos, relativo a la secuencia x , de tal manera que resulte minimizada la suma de los costos de las subsecuencias.

La subsecuencia indizada por $\mathcal{I}(s)$ se denota por $x(s)$. Para establecer el costo asociado a la subsecuencia es de interés conocer la cantidad de veces que cada elemento de \mathcal{X} está en $x(s)$. Esta información se mantiene con un *vector de conteos*, al que denotamos por $c_x(s)$, o $c(s)$, si asumimos que está implícita la secuencia x . O sea, para $a \in \mathcal{X}$

$$c(s)[a] = |\{i \in \mathcal{I}(s) : x_i = a\}|. \quad (2.3)$$

Ejemplo 2.5. Sea $abbabbbc$ la secuencia x procesada. Si el contexto s es ba , entonces $\mathcal{I}(s) = \{3, 6, 8\}$, y $x(s) = bab$. Por lo tanto $c(s)[a] = 1$, $c(s)[b] = 2$ y $c(s)[c] = 0$.

Suponemos que existe una función que a cada vector de conteos le asigna un valor numérico no negativo. Como el vector de conteos está determinado por una subsecuencia, que a su vez es generada por un contexto, podemos extender la definición de esa función al dominio de los contextos. Entonces, existe una función *peso*, w_x , que a cada contexto asigna un valor numérico, real para mayor generalidad, no negativo. Asumiendo que está implícita la secuencia procesada queda definida la función $w: \mathcal{X}^* \rightarrow \mathbb{R}^{\geq 0}$. Se extiende la definición de esta función a un conjunto de contextos: el peso de un conjunto de contextos es la suma de los pesos de los contextos del conjunto. Esto vale en particular para los conjuntos válidos. Sea \mathcal{V} un conjunto válido de contextos relativo a x . Entonces

$$w(\mathcal{V}) = \sum_{s \in \mathcal{V}} w(s). \quad (2.4)$$

Un *conjunto óptimo de contextos* relativo a x es un conjunto válido cuyo peso es menor o igual al peso de cualquier otro conjunto válido.

Árboles

Llamamos *árbol m -ario* a un árbol enraizado en que cada nodo interno tiene exactamente m hijos. Si \mathcal{A} tiene m símbolos, entonces para cada nodo interno se puede etiquetar cada rama saliente del nodo interno con un símbolo diferente de \mathcal{A} . A cada nodo se le puede asociar una cadena de \mathcal{A}^* . Es la cadena que se forma concatenando las etiquetas de las ramas del camino que va desde la raíz hasta el propio nodo. A la raíz corresponde ϵ . Un árbol m -ario es *perfecto* si todas

las hojas están en el mismo nivel. Si la altura de un árbol $|\mathcal{A}|$ -ario perfecto es k , entonces a sus $|\mathcal{A}|^k$ hojas se las puede poner en correspondencia con las cadenas de \mathcal{A}^k .

El procedimiento anterior de etiquetar ramas y nodos se puede generalizar. Las ramas se etiquetan con elementos de \mathcal{A}^+ . Si a es el primer elemento de la etiqueta se dice que se trata de una a -rama. Un \mathcal{A}^+ -árbol es un árbol enraizado cuyas ramas están etiquetadas con elementos de \mathcal{A}^+ y tal que cualquier nodo tiene a lo sumo una a -rama saliente para cada $a \in \mathcal{A}$. En general no son árboles $|\mathcal{A}|$ -arios porque para cualquier nodo interno no necesariamente tiene que existir una a -rama para todo $a \in \mathcal{A}$. También en este caso los nodos se etiquetan con la concatenación de las etiquetas de las ramas del camino desde la raíz al nodo, con ϵ como etiqueta de la raíz. Por simplicidad, con un abuso de lenguaje, identificaremos al nodo con su etiqueta. La rama que va desde el nodo s al nodo t se denota por \vec{st} .

Sea T un \mathcal{A}^+ -árbol y \mathcal{S} el conjunto de cadenas representadas por las hojas de T . Una cadena w es una palabra de \mathcal{T} si $w \preceq s$ para algún $s \in \mathcal{S}$.

Un *árbol de sufijos* T_x para una secuencia x es un \mathcal{A}^+ -árbol cuyas hojas son los sufijos de x . Cualquier subcadena de x es prefijo de un sufijo de x y por lo tanto es una palabra de T_x .

En la Figura 2.1 hay un ejemplo de árbol de sufijos. La existencia de un árbol de sufijos está garantizada si el símbolo x_n no aparece en x_1^{n-1} . De otro modo hay sufijos de x que son prefijos de otros sufijos de x , y por lo tanto no es posible la correspondencia uno a uno entre sufijos y hojas. El problema se evita asumiendo que toda secuencia termina con una marca de finalización, $\#$, que no ocurre en las posiciones previas.

Un *árbol compacto de sufijos* es un árbol de sufijos cuyos nodos internos tienen más de un hijo. En esta tesis árbol de sufijos significa árbol compacto de sufijos.

Se conocen métodos para construcción de árboles de sufijo en tiempo lineal [10, 11, 12, 13].

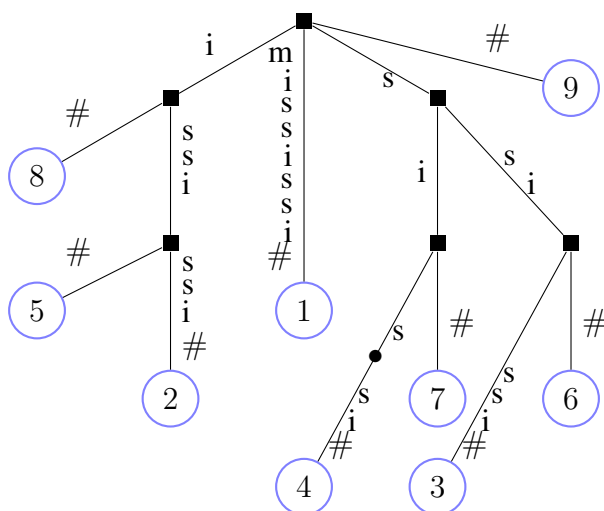


Figura 2.1: Árbol de sufijos para la secuencia 'mississi#'. Hay una hoja para cada sufijo de la secuencia. La concatenación de las etiquetas de las ramas en el camino desde la raíz a la hoja es igual al sufijo. Por ejemplo, el tercer sufijo es *ssissi#* y es la concatenación de *s*, *si* y *ssi#*, el camino desde la raíz a la hoja 3. Cada etiqueta es una subcadena de la secuencia. Sin la marca de finalización '#' el quinto sufijo sería *issi* y no habría una hoja para él, porque sería prefijo de otro sufijo, el segundo, *ississi*. Lo mismo sucede con los sufijos *ssi* y *si*. El octavo sufijo sería *i*, prefijo del segundo y quinto sufijos. Las ramas cuyas etiquetas tienen longitud mayor que 1 pueden ser fragmentadas. Por ejemplo, la rama entrante en la hoja 4 puede fragmentarse en *s* y *si#* creando un nuevo nodo. Pero este nuevo nodo tendría una única rama saliente y el árbol dejaría de ser compacto.

Capítulo 3

Fundamentos sobre conjuntos óptimos de contextos

En este capítulo se considera el problema de encontrar conjuntos óptimos de contextos de una secuencia \mathbf{x} dada cuando la longitud de los contextos está acotada por valores predeterminados. En la Sección 3.1 se trata el caso uni–direccional y en la Sección 3.2 el caso bi–direccional. En el segundo caso se muestra una aplicación en la que los elementos de la secuencia \mathbf{x} pertenecen a \mathcal{A} , y otra en la que pertenecen a $\mathcal{A} \times \mathcal{A}$. En todos los casos se establece con precisión a qué se le llama contexto.

3.1. Contextos uni–direccionales

Aplicaciones típicas de este caso son la compresión y la predicción.

La secuencia procesada $\mathbf{x} = (x_1, x_2, \dots, x_n)$ es simple, es decir, $x_i \in \mathcal{A}$, con $i \in \{1, \dots, n\}$. Los contextos de cada posición son los reversos de las subcadenas de \mathbf{x} que preceden inmediatamente a la posición. Para poder tratar la posición 1 vamos a suponer que hay un elemento x_0 cuyo valor es un símbolo marcador, $\#$, perteneciente a \mathcal{A} , que no ocurre en \mathbf{x} . Entonces, los contextos de la posición i son los prefijos de ${}_{i-1}^0\mathbf{x}$ ($= {}_{i-1}^1\mathbf{x}\#$), que es la cadena de formación de contextos para la posición i , que denotamos por $\mathbf{y}^{(i)}$. Así, para la posición 1 esa cadena es $\#$, y los contextos son ϵ y $\#$. Y para cualquier otra posición i , $\mathbf{y}^{(i)}$ es ${}_{i-1}^0\mathbf{x}$. Estas n cadenas son sufijos de una misma cadena.

Definición 3.1. La cadena de formación de contextos de \mathbf{x}_1^n es $\mathbf{y} = {}_{n-1}^0\mathbf{x}$.

Por lo tanto, para $i \in \{1, \dots, n\}$, la cadena de formación de contextos para la posición i es \mathbf{y}_{n+1-i}^n .

Los contextos son las subcadenas de y . Si el contexto s es una subcadena de y que empieza en la posición j de y , entonces $y_j^{j-1+|s|}$ es un contexto de la posición $n+1-j$ de x . El conjunto de posiciones de x de las que s es contexto es su conjunto índice, $\mathcal{I}(s)$. Nuestro propósito es construir conjuntos válidos de contextos, es decir, aquellos para los cuales los conjuntos índices forman una partición de $\{1, \dots, n\}$.

Ejemplo 3.2. Supongamos que $\mathcal{V} = \{0, 10, 11, \#\}$ es un conjunto de contextos y sea $x = 001010110$, siendo $n = 9$. Entonces $y = 11010100\#$.

x	#	0	0	1	0	1	0	1	1	0
pos		1	2	3	4	5	6	7	8	9
$x(0)$			0	1		1		1		
$x(10)$					0		0		1	
$x(11)$										0
$x(\#)$		0								

y	1	1	0	1	0	1	0	0	#
pos	1	2	3	4	5	6	7	8	9
pos x	9	8	7	6	5	4	3	2	1

Tenemos que $\mathcal{I}(0) = \{2, 3, 5, 7\}$. Estas son las posiciones en x que cumplen que su símbolo precedente es 0. Esto se ve también porque mediante el mapeo $i \mapsto n+1-i$ obtenemos el conjunto $\{8, 7, 5, 3\}$, que son las posiciones de y en donde empieza la subcadena 0. Además $\mathcal{I}(10) = \{4, 6, 8\}$, porque en y esta cadena empieza en las posiciones $\{6, 4, 2\}$. De la misma forma se llega a que $\mathcal{I}(11) = \{9\}$. Por último $\mathcal{I}(\#) = \{1\}$. Se comprueba así que \mathcal{V} es un conjunto válido ya que cada posición pertenece a exactamente uno de los conjuntos índices. La descomposición de x mediante \mathcal{V} genera las subsecuencias $x(0) = (0111)$, $x(10) = (001)$, $x(11) = (0)$ y $x(\#) = (0)$.

Si la secuencia fuese 0010100100 , el conjunto \mathcal{V} seguiría siendo válido aunque en este caso $x(11)$ sería vacía.

Conjunto válido de contextos

Vamos a establecer una condición necesaria y suficiente que deben cumplir los contextos de un conjunto válido de contextos, \mathcal{V} .

Definición 3.3. Se denomina *conjunto libre de prefijos* un conjunto en el que ninguno de sus elementos es prefijo de otro ([14], capítulo 5).

Lema 3.4. *Un conjunto válido de contextos es libre de prefijos.*

Demostración. Supongamos que s y t son contextos de \mathcal{V} y que t es prefijo de s . Para cada posición i de $\mathcal{I}(s)$ el contexto s es prefijo de la cadena de formación de contextos $y^{(i)}$ y, por transitividad, t también es prefijo de $y^{(i)}$. Entonces, cada posición de $\mathcal{I}(s)$ también pertenece a $\mathcal{I}(t)$, por lo que $\mathcal{I}(s) \subseteq \mathcal{I}(t)$. Por lo tanto, \mathcal{V} no es un conjunto de contextos que genere una partición de las posiciones $\{1, \dots, n\}$ por lo que no es válido. \square

De esta manera se llega a una condición necesaria para que un conjunto sea válido, que se conoce como *propiedad de prefijo*, y es que ningún elemento del conjunto sea prefijo de otro:

$$\text{Si } s, t \in \mathcal{V} \text{ con } s \neq t \text{ entonces } s \not\preceq t \text{ y } t \not\preceq s. \quad (3.1)$$

Obviamente, esta no es una condición suficiente para que un conjunto sea válido. Además de que los conjuntos de posiciones sean disjuntos, se debe lograr la *exhaustividad*, esto es, que toda posición pertenezca a uno de ellos. A un conjunto de contextos tal lo llamamos *exhaustivo*. Sea k el máximo de las longitudes de los contextos de \mathcal{V} , $k = \max\{|s| : s \in \mathcal{V}\}$. Denotemos por $\mathcal{P}_k(s)$ el conjunto de extensiones de s de longitud k (versión uni-direccinal de la Expresión (2.2)):

$$\mathcal{P}_k(s) = \{t \in \mathcal{A}^k : s \preceq t\}.$$

Vamos a establecer una condición suficiente para que un conjunto de contextos sea exhaustivo.

En el resto de la tesis se impone como restricción que en cualquier contexto una ocurrencia del símbolo marcador $\#$ sólo puede darse en su última posición. Esta no es una restricción a los contextos reales ya que ese símbolo no ocurre en las secuencias y su presencia es sólo motivada para simplificar estructuras de datos y algoritmos.

Lema 3.5. *Sea \mathcal{V} un conjunto de contextos, tal que $k = \max\{|s| : s \in \mathcal{V}\}$ y*

$$\bigcup_{s \in \mathcal{V}} \mathcal{P}_k(s) = \mathcal{A}^k. \quad (3.2)$$

Entonces se cumple la exhaustividad.

Demostración. Se debe probar que cada posición pertenece al conjunto de índices de un contexto, o sea, que hay un contexto en \mathcal{V} que es prefijo de la cadena de formación de contextos de la posición, esto es, para cada $i \in \{1, \dots, n\} \exists s \in \mathcal{V}, s \preceq y^{(i)}$. Si $y^{(i)}$ tiene longitud k o más sea t su prefijo de longitud k . Entonces, según la Expresión (3.2), existe $s \in \mathcal{V}, s \preceq t$. En otro caso, que la longitud de $y^{(i)}$ sea menor a k implica que su último símbolo es $\#$. Sea $t = y^{(i)}$. Como t

termina en $\#$ ninguna extensión estricta suya puede ser un contexto. Y como sus extensiones de largo k deben tener, otra vez según la Expresión (3.2), un prefijo en \mathcal{V} , se concluye que también en este caso existe $s \in \mathcal{V}$, $s \preceq t$. En ambos casos, por la propiedad transitiva de la relación prefijo, $s \preceq y^{(i)}$. \square

Se deduce el corolario:

Corolario 3.6. *Un conjunto de contextos \mathcal{V} es válido si la relación de extensión aplicada a los elementos s de \mathcal{V} genera una partición del conjunto $\mathcal{P}_k(\epsilon) = \mathcal{A}^k$ en los conjuntos disjuntos $\mathcal{P}_k(s)$, siendo $k = \max\{|s| : s \in \mathcal{V}\}$.*

Esta es una condición suficiente pero no necesaria para que el conjunto \mathcal{V} sea válido. Esto se debe a que alguno de los contextos podría no ocurrir en x . Más adelante, en el Capítulo 4, se tratará esta posibilidad. Por ahora consideramos que un conjunto válido cumple (3.1) y (3.2).

Un conjunto válido de contextos se puede representar con las hojas de un árbol $|\mathcal{A}|$ -ario y también se cumple el recíproco [2]. A estos árboles se les llama, por lo tanto, *árboles de contextos*. Se puede extender la definición de peso a árboles:

Definición 3.7. El *peso de un árbol* es la suma de los pesos de los contextos representados por las hojas.

O sea, el peso de un árbol es igual al del conjunto válido de contextos que representa (Ecuación (2.4)). Entonces, el problema de encontrar un conjunto óptimo de contextos se puede transformar en el problema de hallar un árbol de peso mínimo.

Supongamos que se establece k como una cota superior de las longitudes de los contextos. Entonces un árbol que represente un conjunto válido es un subárbol del árbol $|\mathcal{A}|$ -ario perfecto de altura k .

En la Figura 3.1 vemos un ejemplo de conjunto válido de contextos y árbol de contextos.

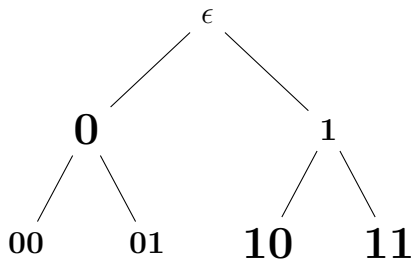


Figura 3.1: Árbol de contextos.

Sea $\mathcal{A} = \{0, 1\}$. Vemos un árbol perfecto de contextos para $k = 2$. El conjunto $\{0, 10\}$ no es exhaustivo para cualquier secuencia, porque 11 no es extensión de ningún elemento del conjunto. El conjunto $\{0, 10, 11, 01\}$ no es disjunto porque 01 es extensión de 0 y de 01. Un conjunto válido para cualquier secuencia es, por ejemplo, $\{0, 10, 11\}$ que está marcado en negrita en la figura.

Para encontrar el árbol óptimo para cada cadena \mathbf{x} se aplica un algoritmo que empieza en las hojas del árbol perfecto. Se define el *peso óptimo* de cada contexto \mathbf{s} , $w_{opt}(\mathbf{s})$. Si el nodo \mathbf{s} es una hoja se calcula el peso con la función $w(\mathbf{s})$, y $w_{opt}(\mathbf{s}) = w(\mathbf{s})$. Sea \mathbf{s} un nodo interno del que ya fueron procesados sus hijos. Se calcula su peso óptimo como el mínimo entre el peso del contexto y la suma de los pesos óptimos de sus hijos: $w_{opt}(\mathbf{s}) = \min \{w(\mathbf{s}), \sum_{a \in \mathcal{A}} w_{opt}(\mathbf{s}a)\}$. Se construye el árbol óptimo de ese nodo: si el mínimo es el primero de los valores, el árbol es podado y consiste en ese nodo. En el otro caso el árbol es el nodo junto con los árboles óptimos que tienen raíz en los hijos del nodo considerado. En el ejemplo de la Figura 3.1 esto último es lo que pasa en los nodos ϵ y 1 si el conjunto representado es el óptimo. En el caso en que los dos valores sean iguales se elige de manera determinística una de las opciones, por ejemplo, el árbol podado.

Un árbol que corresponde a un conjunto válido describe de manera recursiva la partición del conjunto $\mathcal{P}_k(\epsilon)$. Si el árbol no es sólo una hoja, en la raíz se obtienen los $|\mathcal{A}|$ conjuntos $\mathcal{P}_k(a)$ para cada $a \in \mathcal{A}$. Recursivamente, en cada nodo interno \mathbf{s} se establece la partición de $\mathcal{P}_k(\mathbf{s})$ en los conjuntos $\mathcal{P}_k(\mathbf{s}a)$.

3.2. Contextos bi–direccionales

Se ha generalizado el modelado por contextos, considerando contextos compuestos por múltiples componentes profundizando en el caso *bi–direccional* [5, 6, 7]. Para cada posición i de \mathbf{x} se dispone de dos cadenas formadoras de contexto, $\mathbf{y}_\ell^{(i)}$ y $\mathbf{y}_r^{(i)}$, llamadas “izquierda” y “derecha” respectivamente. Los *contextos bi–direccionales*, o *bi–contextos* son pares ordenados de contextos, $(\mathbf{s}_\ell, \mathbf{s}_r)$.

En la generalización más directa del caso uni–direccional, la secuencia \mathbf{x} consta de dos pistas, $\mathbf{x} = (\mathbf{x}_\ell, \mathbf{x}_r)$, de tal modo que $\mathbf{x}_i = ([\mathbf{x}_\ell]_i, [\mathbf{x}_r]_i)$. Tal como en el caso uni–direccional se supone que hay un elemento en la posición 0, $\mathbf{x}_0 = (\#, \#)$, siendo $\#$ un elemento de \mathcal{A} que no ocurre en la cadena de entrada. Las cadenas de formación de contextos son $\mathbf{y}_\ell = {}_{n-1}^0[\mathbf{x}_\ell]$, y $\mathbf{y}_r = {}_{n-1}^0[\mathbf{x}_r]$.

En la aplicación eliminación de ruido, los elementos de \mathbf{x} son símbolos del alfabeto. Para cada posición i , los contextos izquierdos son los reversos de las subcadenas que preceden la posición i , mientras que los contextos derechos son las subcadenas que empiezan en $i + 1$. También aquí se asume que el elemento $\#$ de \mathcal{A} no ocurre en la cadena y se supone $\mathbf{x}_0 = \#$ y $\mathbf{x}_{n+1} = \#$. Entonces \mathbf{y}_ℓ es ${}_{n-1}^0\mathbf{x} = {}_{n-1}^1\mathbf{x}\#$, como antes, pero \mathbf{y}_r es $\mathbf{x}_2^{n+1} = \mathbf{x}_2^n\#$ y los contextos derechos de la posición i , con $1 \leq i \leq n$, son los prefijos de $[\mathbf{y}_r]_i^n$.

El contexto $(\mathbf{s}_\ell, \mathbf{s}_r)$ ocurre en la posición i si es prefijo de $(\mathbf{y}_\ell^{(i)}, \mathbf{y}_r^{(i)})$. Las posiciones en que ocurre se denotan por $\mathcal{I}(\mathbf{s}_\ell, \mathbf{s}_r)$ y la subsecuencia indizada por esas posiciones es $\mathbf{x}(\mathbf{s}_\ell, \mathbf{s}_r)$.

Se mantiene que un conjunto de contextos \mathcal{V} es válido si genera una partición de las posiciones $\{1, \dots, n\}$. La condición para que esto se cumpla es la generalización de las propiedades del caso uni–direccional. Esto es:

- (1) \mathcal{V} debe ser un conjunto libre de prefijos, lo que implica que los conjuntos de extensiones sean disjuntos por pares: $\mathcal{P}(s_\ell, s_r) \cap \mathcal{P}(t_\ell, t_r) = \emptyset$ si $(s_\ell, s_r) \neq (t_\ell, t_r)$, con $(s_\ell, s_r), (t_\ell, t_r) \in \mathcal{V}$;
- (2) \mathcal{V} debe ser exhaustivo, esto es, si k_ℓ (resp. k_r) es la mayor longitud de todos los componentes izquierdos (resp. derechos) de los contextos de \mathcal{V} , entonces cada contexto con longitud de componente izquierdo (resp. derecho) igual a k_ℓ (resp. k_r) es extensión de algún elemento de \mathcal{V} : $\mathcal{P}_{(k_\ell, k_r)}(\epsilon, \epsilon) \subseteq \bigcup_{(s_\ell, s_r) \in \mathcal{V}} \mathcal{P}(s_\ell, s_r)$. En el resto de esta sección k_ℓ, k_r tienen el significado aquí definido.

Se ha demostrado [5, 6, 7] que, tal como en el caso uni–direccional, un conjunto válido de contextos bi–direccionales puede representarse mediante las hojas de un árbol cuyos nodos pertenecen ahora a $\mathcal{A}^* \times \mathcal{A}^*$, con (ϵ, ϵ) como raíz. También se cumple el recíproco, un árbol de contextos bi–direccionales representa un conjunto válido de contextos bi–direccionales. Sin embargo, la estructura de este árbol es diferente al del caso uni–direccional. Los *árboles de contextos bi–direccionales* cumplen que cada nodo interno tiene como hijos o bien todas sus extensiones mínimas por izquierda o bien todas sus extensiones mínimas por derecha. O sea, si (s_ℓ, s_r) no es una hoja, sus hijos son o bien $\{(s_\ell a, s_r) : a \in \mathcal{A}\}$, o bien $\{(s_\ell, s_r a) : a \in \mathcal{A}\}$.

A diferencia del caso uni–direccional, un conjunto válido puede ser representado por más de un árbol.

Por ejemplo, el conjunto $\{(0, 0), (1, 01), (10, 1), (01, 1), (1, 1)\}$, con $\mathcal{A} = \{0, 1\}$, no es exhaustivo porque ningún contexto es prefijo de $(10, 00)$. Y no es libre de prefijos porque $(10, 1)$ es extensión de $(1, 1)$ y $(10, 1)$.

El conjunto $\{(0, 0), (1, 00), (1, 01), (00, 1), (01, 1), (1, 1)\}$ es válido. Y puede ser representado por dos árboles, como se muestra en la Figura 3.2.

Como en el caso uni–direccional un conjunto válido puede describirse con una secuencia de particiones de las extensiones pertenecientes a $\mathcal{P}_{(k_\ell, k_r)}(\epsilon, \epsilon)$ de los nodos del árbol que representa al conjunto. Ahora hay tres posibles particiones. La “no partición”, que es lo que se hace en las hojas del árbol, y las particiones izquierda y derecha.

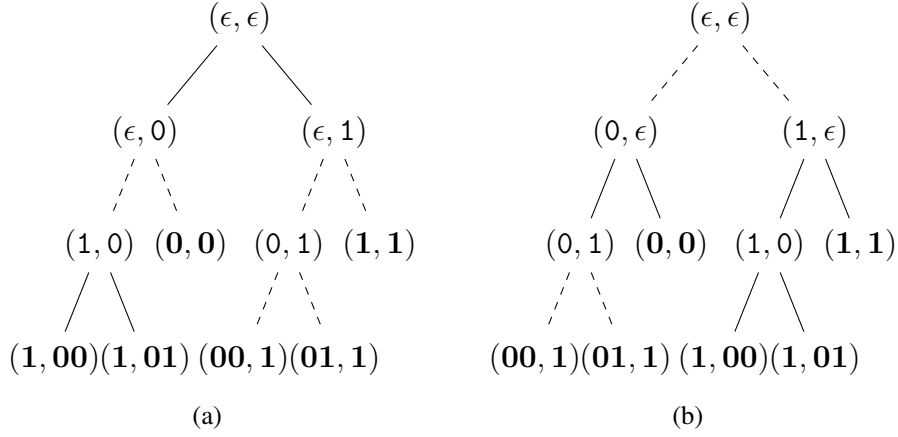


Figura 3.2: Dos árboles de contextos bi–direccionales que representan el mismo conjunto de contextos, que son las hojas de los árboles, iguales en ambos. En trazo continuo se muestran las particiones derechas y en trazo discontinuo las particiones izquierdas.

Podado bi–direccional

Tal como en el caso uni–direccional las subsecuencias $\mathbf{x}(s_\ell, s_r)$ determinadas por los contextos (s_ℓ, s_r) inducen un costo por lo que a los contextos se les asigna un peso $w(s_\ell, s_r)$. El peso del conjunto válido de contextos es la suma de los pesos de los contextos y queremos encontrar un conjunto óptimo, esto es, de peso mínimo.

Con $\mathcal{T}(s_\ell, s_r)$ representamos los subárboles con raíz en (s_ℓ, s_r) de los árboles de contextos bi–direccionales. Las hojas de un árbol $T \in \mathcal{T}(s_\ell, s_r)$ son contextos que forman un subconjunto de un conjunto válido. Los elementos de este subconjunto son los contextos del conjunto válido que son extensión de (s_ℓ, s_r) . Esto motiva la siguiente definición.

Definición 3.8. El conjunto de contextos \mathcal{V} es un *conjunto válido relativo a* (s_ℓ, s_r) si el conjunto $\{\mathcal{P}_{(k_\ell, k_r)}(t_\ell, t_r) : (t_\ell, t_r) \in \mathcal{V}\}$ es una partición de $\mathcal{P}_{(k_\ell, k_r)}(s_\ell, s_r)$, donde k_ℓ (resp. k_r) es el máximo de las longitudes de los componentes izquierdos (resp. derechos) de los elementos de \mathcal{V} .

Es claro que debe cumplirse $(s_\ell, s_r) \preceq \mathcal{V}$.

También podemos definir un *conjunto óptimo relativo a* (s_ℓ, s_r) como uno que tenga peso mínimo entre los conjuntos válidos relativos a (s_ℓ, s_r) .

Obviamente, un conjunto válido y un conjunto óptimo como fue definido antes es un conjunto válido relativo a (ϵ, ϵ) y un conjunto óptimo relativo a (ϵ, ϵ) respectivamente. Esto sugiere un método para encontrar un conjunto óptimo (o

sea, un conjunto óptimo relativo a la raíz de un árbol) si se puede demostrar que se cumple el principio de optimalidad. Se empieza con los contextos más largos, para cada uno de los cuales el conjunto óptimo relativo a él es el propio contexto. Para cualquier otro contexto se va a conocer un conjunto óptimo relativo a cada una de sus extensiones mínimas, tanto izquierdas como derechas. Haciendo valer el principio de optimalidad se obtiene, mediante un algoritmo de programación dinámica [15], el conjunto óptimo relativo al contexto procesado.

Con $\mathcal{T}_{opt}(s_\ell, s_r)$ se denota al subconjunto de $\mathcal{T}(s_\ell, s_r)$ cuyo peso es mínimo, siendo el peso de un árbol igual a la suma de los pesos de sus hojas, esto es, el peso del conjunto válido que representa. Las hojas de un árbol de $\mathcal{T}_{opt}(s_\ell, s_r)$ forman un conjunto óptimo relativo a (s_ℓ, s_r) . Como ya fue notado, $\mathcal{T}_{opt}(s_\ell, s_r)$ puede tener más de un elemento, aún en el caso de que haya un único conjunto válido de peso mínimo. Cualquiera de los árboles pertenecientes a $\mathcal{T}_{opt}(s_\ell, s_r)$ puede ser denotado por $T_{opt}(s_\ell, s_r)$. Al peso de los árboles de $\mathcal{T}_{opt}(s_\ell, s_r)$ lo denotamos por $w_{opt}(s_\ell, s_r)$.

Se ha demostrado [16, 6, 7] que se cumple el principio de optimalidad que nos permite establecer el peso óptimo de un contexto conociendo el de sus extensiones mínimas:

$$w_{opt}(s_\ell, s_r) = \min \left\{ w(s_\ell, s_r), \sum_{a \in \mathcal{A}} w_{opt}(s_\ell a, s_r), \sum_{a \in \mathcal{A}} w_{opt}(s_\ell, s_r a) \right\}. \quad (3.3)$$

Si en el segundo término $|s_\ell a| > k_\ell$, o sea, si el contexto no cumple la restricción establecida, entonces se considera que $w_{opt}(s_\ell a, s_r) = \infty$. Lo mismo sucede en el tercer término si $|s_r a| > k_r$.

Se debe establecer una regla para elegir el mínimo en caso de que más de uno de estos valores lo sea. Y entonces se puede definir la función `split` para cada contexto bi-direccional que refleja cual es la opción que permite obtener el peso óptimo. El valor de esta función es $\{ 'N', 'L', 'R' \}$ según que el óptimo se alcance respectivamente mediante la “no partición”, la partición izquierda o la partición derecha.

Entonces, como en el caso uni-direccional se puede establecer un algoritmo que mediante programación dinámica encuentra un árbol óptimo. Este es el Algoritmo 1.

Cuando $j = 0$ se cumple que $(s_\ell, s_r) = (\epsilon, \epsilon)$ y por lo tanto se obtiene el árbol óptimo. Los árboles se están describiendo mediante sus nodos, sin mencionar las ramas que quedan implícitamente definidas. Al tener reglas de desempate en cada conjunto $\mathcal{T}_{opt}(s_\ell, s_r)$ hay un árbol que el algoritmo obtiene de manera determinística que se denota por $T_{opt}(s_\ell, s_r)$.

El peso de un contexto es una función del vector de conteos del contexto. En el primer bucle del Algoritmo 1 se calcula el peso de los contextos de longitud

Algoritmo 1 Árbol óptimo de contextos bi–direccionales.

Para cada $(\mathbf{s}_\ell, \mathbf{s}_r) \in \mathcal{A}^{k_\ell} \times \mathcal{A}^{k_r}$

calcular $w(\mathbf{s}_\ell, \mathbf{s}_r)$, y asignar $w_{opt}(\mathbf{s}_\ell, \mathbf{s}_r) = w(\mathbf{s}_\ell, \mathbf{s}_r)$

$T_{opt}(\mathbf{s}_\ell, \mathbf{s}_r) = \{(\mathbf{s}_\ell, \mathbf{s}_r)\}$

Para $j = k_\ell + k_r - 1$ **hasta** 0

Para cada $(\mathbf{s}_\ell, \mathbf{s}_r)$ tal que $|\mathbf{s}_\ell| + |\mathbf{s}_r| = j, |\mathbf{s}_\ell| \leq k_\ell, |\mathbf{s}_r| \leq k_r$

calcular $w(\mathbf{s}_\ell, \mathbf{s}_r)$ y, mediante la ecuación (3.3), $w_{opt}(\mathbf{s}_\ell, \mathbf{s}_r)$

obtener $\text{split}(\mathbf{s}_\ell, \mathbf{s}_r)$

$T_{opt}(\mathbf{s}_\ell, \mathbf{s}_r) =$

$$\begin{cases} \{(\mathbf{s}_\ell, \mathbf{s}_r)\} & \text{si } \text{split}(\mathbf{s}_\ell, \mathbf{s}_r) = \text{'N'}, \\ \{(\mathbf{s}_\ell, \mathbf{s}_r)\} \cup \bigcup_{a \in \mathcal{A}} T_{opt}(\mathbf{s}_\ell a, \mathbf{s}_r) & \text{si } \text{split}(\mathbf{s}_\ell, \mathbf{s}_r) = \text{'L'}, \\ \{(\mathbf{s}_\ell, \mathbf{s}_r)\} \cup \bigcup_{a \in \mathcal{A}} T_{opt}(\mathbf{s}_\ell, \mathbf{s}_r a) & \text{si } \text{split}(\mathbf{s}_\ell, \mathbf{s}_r) = \text{'R'}. \end{cases}$$

máxima, por lo cual debe obtenerse el vector de conteos mediante inspección de \mathbf{x} . Pero para el resto de los contextos no es necesario porque el vector se obtiene a partir de los contextos más largos.

$$\mathbf{c}(\mathbf{s}_\ell, \mathbf{s}_r) = \sum_{a \in \mathcal{A}} \mathbf{c}(\mathbf{s}_\ell a, \mathbf{s}_r) = \sum_{a \in \mathcal{A}} \mathbf{c}(\mathbf{s}_\ell, \mathbf{s}_r a). \quad (3.4)$$

Hay que notar una diferencia con el caso uni–direccional. En aquel caso el algoritmo que obtiene el árbol óptimo recorre desde abajo hacia arriba, desde los contextos largos hacia los cortos, un árbol, de hecho un árbol perfecto. En la generalización bi–direccional, la estructura implícita que se está recorriendo es una generalización del árbol perfecto. Lo llamamos *grafo de contextos bi–direccionales*. Cada nodo interno tiene dos conjuntos de hijos. Si los dos componentes de un nodo son diferentes de ϵ , entonces el nodo tiene dos ramas entrantes en él. El nodo $(\mathbf{s}_\ell a, \mathbf{s}_r b)$ desciende de $(\mathbf{s}_\ell a, \mathbf{s}_r)$ y de $(\mathbf{s}_\ell, \mathbf{s}_r b)$, o sea, tiene dos “padres”. Respecto al primero está incluido en la partición derecha, y respecto al segundo en la partición izquierda.

Advertimos otra diferencia con el caso uni–direccional que, como se verá en el Capítulo 4, determina la complejidad de obtener un conjunto óptimo. En una partición de un nodo interno, digamos la izquierda, no están incluidas todas las extensiones estrictas, sino sólo las extensiones estrictas en el componente izquierdo. Por ejemplo en la Figura 3.2(b), en la raíz se elige la partición izquierda y, como consecuencia, el contexto $(\epsilon, 0)$ no está incluido en la partición. Sin embargo las extensiones $\mathcal{P}_{(k_\ell, k_r)}(\epsilon, 0)$ de $(\epsilon, 0)$, con $k_\ell, k_r \geq 2$, si están incluidos en la

partición. Esto es porque ese conjunto de extensiones son también extensión de otros contextos.

Se debe destacar la importancia de poder representar un conjunto válido mediante un árbol porque existen ejemplos para 3 direcciones en que eso no es cierto [7] y se ha demostrado que ese no es un caso aislado para multi-direccionalidad mayor a 2 [6].

Lo demostrado en el Corolario 3.6 y su equivalente en el marco bi-direccional corresponde a conjuntos válidos para cualquier secuencia. En lo que sigue el concepto conjunto válido es aplicado para cada secuencia específica y se excluyen del conjunto aquellos contextos que no ocurren en la secuencia.

Capítulo 4

Contextos Representativos

En este capítulo nos planteamos el problema de eliminar la restricción que en el primer bucle del Algoritmo 1 se impone a la longitud de los contextos. Esto en realidad significa que el límite sea n . El problema es que el algoritmo, aún si se simplifica procesando sólo los contextos que ocurren, en el caso bi-direccional es de complejidad cúbica, ya que para cada una de las n posiciones hay $\Theta(n^2)$ contextos.

En el caso uni-direccional se puede resolver el problema, y la complejidad cuadrática de la solución simplista se reduce a lineal utilizando árboles de sufijos [4]. Tal como se describirá en el Capítulo 6, aquí utilizaremos esa misma herramienta y además algoritmos para calcular el prefijo común más largo de los contextos representados en los árboles.

Veremos que los contextos se pueden separar en clases de tal modo que todos los elementos de la misma clase tengan las mismas propiedades, determinadas por ocurrir en exactamente las mismas posiciones de x . Entonces, en lugar de procesar todos los contextos que ocurren se procesa sólo uno por clase, el representativo.

Vamos a caracterizar las clases y los contextos representativos. Hay un equivalente a los contextos de longitud máxima con el cual comienza el Algoritmo 1. Son aquellos que no pueden extenderse porque alcanzan el final de la cadena de formación de contextos. Los contextos de la misma clase son aquellos que son prefijos de los mismos contextos de longitud máxima. Estos contextos de longitud máxima son la generalización de las hojas de los árboles de sufijos. La generalización de los nodos internos se obtiene como un subconjunto del producto cartesiano de los nodos internos de los árboles de sufijos que corresponden a las dos cadenas de formación de contextos. Demostraremos que en la búsqueda de un conjunto óptimo de contextos para la secuencia x dada alcanza con considerar sólo los contextos representativos.

4.1. Caso uni–direccional

Para el caso uni–direccional habíamos visto que las cadenas de formación de contextos de todas las posiciones son sufijos de una misma cadena: el reverso de \mathbf{x} con el agregado de la marca de finalización $\#$, obteniendo ${}_{n-1}^0\mathbf{x}(= {}_{n-1}^1\mathbf{x}\#)$. A esa cadena, la denotamos por \mathbf{y} , y la llamamos la *cadena de formación de contextos* de \mathbf{x} . Por lo tanto, los contextos de \mathbf{x}_i son los prefijos de \mathbf{y}_{n+1-i}^n . Decimos que un contexto s ocurre en la posición i si $s \preceq \mathbf{y}_{n+1-i}^n$.

Definición 4.1. Para cada i , $1 \leq i \leq n$, el contexto de máxima longitud que ocurre en i es llamado *contexto sufijo*.

El agregado de la marca de finalización es para poder representar la cadena \mathbf{y} con un árbol de sufijos, $T_{\mathbf{y}}$. Los contextos sufijos terminan en $\#$ y, como ya se vio, están en correspondencia uno a uno con las hojas de $T_{\mathbf{y}}$.

Un contexto puede ocurrir en diferentes posiciones, y por lo tanto puede ser prefijo de varios contextos sufijos. Recordemos que llamamos conjunto índice de s y lo denotamos con $\mathcal{I}(s)$ al conjunto de posiciones donde ocurre s .

Definición 4.2. El *conjunto de contextos sufijo* de \mathbf{y} determinado por el contexto s , denotado por $\mathcal{S}(s)$, es el conjunto de sufijos de \mathbf{y} que extienden s . Esto es

$$\mathcal{S}(s) = \{\mathbf{y}_{n+1-i}^n, i \in \mathcal{I}(s)\}.$$

Entonces s ocurre si y sólo si $\mathcal{S}(s)$ no es vacío.

Como $s \preceq \mathcal{S}(s)$ y la relación prefijo es transitiva, se tiene que

$$\mathbf{t} \preceq \mathbf{s} \Rightarrow \mathbf{t} \preceq \mathcal{S}(s) \Rightarrow \mathcal{S}(s) \subseteq \mathcal{S}(\mathbf{t}). \quad (4.1)$$

El recíproco no se cumple en general porque puede haber varios contextos con el mismo conjunto sufijo. Esto es lo que permitirá reducir la complejidad porque sólo hará falta procesar uno de ellos.

Ejemplo 4.3. Sea $\mathbf{x} = 1110101$, y por lo tanto $\mathbf{y} = 010111\#$. El contexto 011 ocurre en la posición 5 (ya que $011 = x_4x_3x_2$), por lo que $\mathcal{S}(011) = \{0111\#$. Sea $s = 0$ un prefijo de 011 . Entonces, $\mathcal{I}(s)$ es $\{5, 7\}$, de donde $\mathcal{S}(s) = \{0111\#, 010111\# \} \supseteq \mathcal{S}(011)$.

Para ver que no se cumple el recíproco de (4.1) sea $\mathbf{t} = 01$. Tenemos que \mathbf{t} ocurre en las mismas posiciones que s . Por lo tanto se cumple $\mathcal{S}(s) \subseteq \mathcal{S}(\mathbf{t})$, aunque $\mathbf{t} \not\preceq s$.

Como hay una correspondencia entre las posiciones y los sufijos, es inmediato el siguiente teorema.

Teorema 4.4. *Se cumple $\mathcal{S}(s) = \mathcal{S}(t)$ si y sólo si $\mathcal{I}(s) = \mathcal{I}(t)$.*

La relación de igualdad entre conjuntos de sufijos define clases de equivalencia.

Definición 4.5. *La partición por igualdad de conjuntos de sufijos de la secuencia x , es la partición en clases del conjunto de contextos que ocurren, de tal manera que todos los elementos de una clase tienen el mismo conjunto de sufijos.*

Denotamos la partición por Π , y una clase genérica por \mathcal{C} . La clase a la cual pertenece el contexto s se denota por $\mathcal{C}(s)$. Además, denotamos por $\mathcal{S}(\mathcal{C})$ al conjunto de contextos sufijos de los elementos de la clase \mathcal{C} .

Ejemplo 4.6. Según lo visto en el Ejemplo 4.3, 0 y 01 pertenecen a una clase \mathcal{C} , tal que $\mathcal{S}(\mathcal{C}) = \{0111\#, 010111\#\}$. Ya sabíamos que $\mathcal{S}(011)$ está contenido de manera estricta en $\mathcal{S}(\mathcal{C})$, por lo que 011 no está en \mathcal{C} . Se puede ver que $\mathcal{S}(010) = \{010111\#\}$, y por lo tanto, tampoco 010 pertenece a \mathcal{C} . Entonces ninguna extensión estricta de 01 puede pertenecer a \mathcal{C} (ya que no pertenecen ni 010 ni 011 y según la Expresión (4.1) los conjuntos de sufijos de las extensiones de estos contextos están contenidos en $\mathcal{S}(010)$ o $\mathcal{S}(011)$). Por otro lado 00 no ocurre en x . Finalmente $\mathcal{S}(\epsilon)$ es un conjunto que contiene siete elementos, que son todos los sufijos de y , por lo que ϵ no pertenece a \mathcal{C} . Entonces en la clase \mathcal{C} no hay otros contextos, esto es, $\mathcal{C} = \{0, 01\}$.

Una clase queda determinada por un conjunto de contextos sufijos. Pero cualquier conjunto de contextos sufijos no determina una clase.

Teorema 4.7. *Dada una clase \mathcal{C} de Π todos los elementos de $\mathcal{S}(\mathcal{C})$ deben tener un prefijo común, que no debe ser prefijo de ningún otro contexto sufijo.*

Demostración. La primera parte de la afirmación es evidente ya que $\epsilon \preceq \mathcal{S}(\mathcal{C})$. Sea s el más largo de los contextos tal que $s \preceq \mathcal{S}(\mathcal{C})$. Supongamos que además existiera un contexto sufijo t , tal que $t \notin \mathcal{S}(\mathcal{C})$ y $s \preceq t$.

Por ser s el más largo de los prefijos comunes de $\mathcal{S}(\mathcal{C})$, para toda extensión estricta de s , sa , $a \in \mathcal{A}$, existe algún elemento de $\mathcal{S}(\mathcal{C})$, sbv , $b \in \mathcal{A}$, $b \neq a$, del cual no es prefijo. Por lo tanto sa no pertenece a \mathcal{C} . Entonces, los elementos de \mathcal{C} deben ser prefijos de s . Pero todos los prefijos de s , con este incluido, son también prefijos de t . Esto contradice la hipótesis de que $t \notin \mathcal{S}(\mathcal{C})$. \square

Ejemplo 4.8. En la Figura 4.1 se muestra el árbol de sufijos T_y correspondiente a la cadena de formación de contextos $y = 010111\#$ del Ejemplo 4.3. El conjunto de sufijos representado por las hojas 1 y 3 determinan la clase $\{0, 01\}$.

La hoja 3 determina la clase $\{011, 0111, 0111\#\}$.

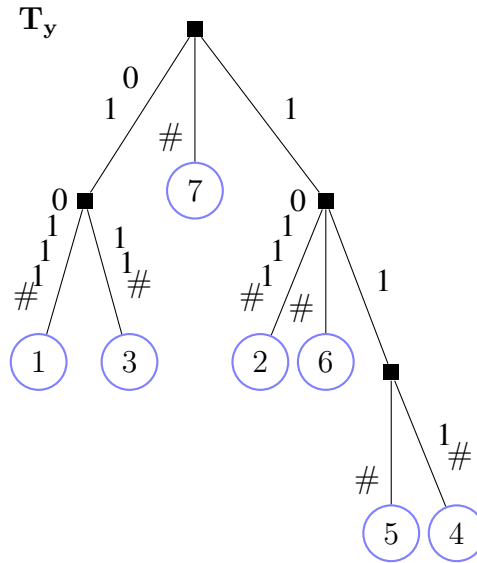


Figura 4.1: Árbol de la cadena de formación de contextos de $y = 010111\#$.

Ahora veamos que los sufijos correspondientes a las hojas 2, 4 y 5 no determinan una clase. Los únicos contextos que son prefijos de estos tres sufijos son ϵ y 1. Entonces estos son los únicos elementos posibles de la hipotética clase. Pero 1 es también prefijo del sufijo correspondiente a la hoja 6. Y ϵ , además es también prefijo de 1, 3, 6 y 7. O sea que ningún contexto es prefijo de 2, 4 y 5 pero de ningún otro sufijo. Por otra parte 2, 4, 5 y 6 determinan la clase $\{1\}$.

Vamos a caracterizar las clases y definir los contextos representativos.

Cada rama \vec{ts} de T_y define una cadena $\mathbf{u} = u_1u_2 \dots u_j, j \geq 1$, tal que $\mathbf{s} = \mathbf{t}\mathbf{u}$. Por lo tanto \vec{ts} representa al conjunto de contextos $\mathcal{C} = \{\mathbf{t}\mathbf{u}_1^i, 1 \leq i \leq j\}$. Se debe notar que \mathbf{s} ($= \mathbf{t}\mathbf{u}_1^j$) pertenece a este conjunto.

Teorema 4.9. *Sea s un nodo de T_y , distinto de la raíz, y t su nodo padre. Los contextos representados por \vec{ts} son una clase.*

Demostración. Las hojas (sufijos) que descenden de esos contextos son las mismas. Falta probar que ningún otro contexto es prefijo de exactamente esas hojas. Otros contextos que también son prefijos de esas hojas son t y sus prefijos. Pero estos contextos son además prefijos de otras hojas. Si s es una hoja el teorema queda probado. En otro caso cualquier par de extensiones mínimas de s son disjuntas y por lo tanto ninguna de ellas ni sus extensiones son prefijos de las mismas hojas. \square

Corolario 4.10. *Hay una clase de Π para cada nodo de T_y .*

Demostración. Cada nodo distinto de la raíz identifica la rama que entra en él y por lo tanto a una clase. El nodo raíz es otra clase ya que es el único contexto del cual descenden todas las hojas. \square

Vemos que el conjunto \mathcal{C} acepta un orden total en que cada contexto es prefijo del siguiente. Cualquier contexto de \mathcal{C} puede representar la clase. Se puede elegir tu_1 [4], que es prefijo (el prefijo común más largo) de todos los elementos de la clase. Para nuestra generalización al caso bi–direccional tendremos que elegir otro porque no habrá, en general, un contexto que sea prefijo de todos los otros de la clase. Anticipamos que el que elegimos es el equivalente a s , que es la extensión común más corta (concepto que se presentará en la Definición 4.15) de todos los elementos de la clase.

Definición 4.11. El *contexto representativo* de la clase \mathcal{C} , denotado por $R(\mathcal{C})$ es el contexto s , que es la extensión de todos los contextos de la clase. También, por extensión, definimos el representativo del contexto u y lo denotamos por $R(u)$ como el representativo de la clase a la que pertenece u .

4.2. Generalización Bi–direccional

En el caso bi–direccional asumimos que, al igual que en el uni–direccional, todas las cadenas $y_\ell^{(i)}$ y $y_r^{(i)}$ son sufijos respectivamente de y_ℓ y y_r llamadas *cadena de formación de contextos izquierda y derecha* respectivamente. De esta manera todos los contextos izquierdos están representados en un árbol de sufijos, T_{y_ℓ} y todos los contextos derechos están representados en T_{y_r} .

Definición 4.12. Los *contextos bi–sufijo* son los n contextos $(y_\ell^{(i)}, y_r^{(i)})$, con $i \in \{1, \dots, n\}$. Cada uno de estos es el contexto más largo para una de las posiciones.

Recordamos que el contexto bi–direccional (s_ℓ, s_r) ocurre en i si es prefijo de $(y_\ell^{(i)}, y_r^{(i)})$. Definimos $\mathcal{I}(s_\ell, s_r)$ como en el caso uni–direccional. Tal como en el caso uni–direccional se define el *conjunto de sufijos* $\mathcal{S}(s_\ell, s_r)$, esto es, el conjunto de contextos bi–sufijo que son extensión de (s_ℓ, s_r) . Definimos asimismo la partición por igualdad de conjuntos de sufijos Π del conjunto de todos los contextos bi–direccionales que ocurren en x . De esta manera si (s_ℓ, s_r) y (t_ℓ, t_r) pertenecen a la misma clase \mathcal{C} , entonces $\mathcal{S}(s_\ell, s_r) = \mathcal{S}(t_\ell, t_r) = \mathcal{S}(\mathcal{C})$. La clase a la que pertenece (s_ℓ, s_r) la denotamos $\mathcal{C}(s_\ell, s_r)$. Además \mathcal{S} denota al conjunto de todos los bi–sufijos. Cualquier contexto bi–direccional que ocurre es prefijo de un elemento de \mathcal{S} . También se cumple la generalización de la Expresión (4.1):

$$(t_\ell, t_r) \preceq (s_\ell, s_r) \Rightarrow (t_\ell, t_r) \preceq \mathcal{S}(s_\ell, s_r) \Rightarrow \mathcal{S}(s_\ell, s_r) \subseteq \mathcal{S}(t_\ell, t_r). \quad (4.2)$$

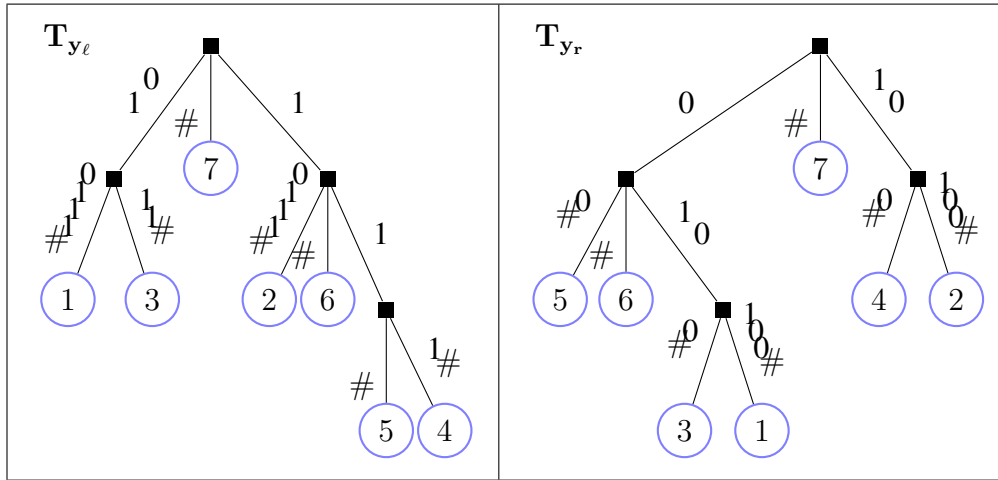


Figura 4.2: Árboles de las cadenas de formación de contextos para la cadena $x = \begin{pmatrix} 1110101 \\ 0010101 \end{pmatrix}$.

Notemos que un bi-sufijo está determinado por dos contextos que son sufijos en sus respectivas cadenas de formación de contextos, y por lo tanto hojas en sus respectivos árboles de sufijos. Por lo tanto la relación bi-sufijos induce una correspondencia uno a uno entre las n hojas de T_{y_ℓ} y las n hojas de T_{y_r} .

Ejemplo 4.13. Sea $x = \begin{pmatrix} 1110101 \\ 0010101 \end{pmatrix}$, de donde $\begin{pmatrix} y_\ell \\ y_r \end{pmatrix} = \begin{pmatrix} 010111\# \\ 010100\# \end{pmatrix}$. Los contextos $(01, 0)$ y $(0, 010)$ ocurren ambos en las mismas posiciones, 5 y 7, de x y por lo tanto pertenecen a la misma clase \mathcal{C} .

En la Figura 4.2 se muestran los árboles de sufijos de las cadenas formadoras de contextos.

4.3. Contextos Representativos

Mientras que en el caso uni-direccional, como ya se mencionó, los contextos que pertenecen a la misma clase \mathcal{C} satisfacen una relación de orden total, esto *no* sucede en el caso bi-direccional. Más aún, no existe un contexto bi-direccional que sea prefijo de todos los otros de la clase. Sin embargo, como veremos en el Lema 4.17, cada clase contiene un elemento que es una extensión común de todos

los contextos bi–direccionales de la clase. Antes de establecer el lema veremos un ejemplo y agregaremos una definición.

Ejemplo 4.14. En el Ejemplo 4.13, además de $(01, 0)$ y $(0, 010)$, también pertenecen a la clase \mathcal{C} los contextos bi–direccionales $(0, \epsilon)$, $(01, \epsilon)$, $(0, 0)$, $(\epsilon, 01)$, $(\epsilon, 010)$, $(0, 01)$, $(01, 01)$, $(01, 010)$. La presencia de $(0, \epsilon)$ y $(\epsilon, 01)$ evidencia que no hay un contexto bi–direccional, distinto de (ϵ, ϵ) , que sea prefijo de todos los demás. Por otra parte, el contexto bi–direccional $(01, 010)$ es una extensión de todos los contextos bi–direccionales de \mathcal{C} .

Definimos la *extensión común más corta* de dos contextos uni–direccionales s y t no disjuntos como el más largo de ellos. En un conjunto de contextos no-disjuntos por pares la extensión común más corta es el más largo de los contextos. La motivación de esta definición es su generalización.

Definición 4.15. Sean (s_ℓ, s_r) y (t_ℓ, t_r) contextos bi–direccionales no–disjuntos. La *extensión común más corta* de (s_ℓ, s_r) y (t_ℓ, t_r) es el contexto bi–direccional cuyo componente izquierdo (resp. derecho) es la extensión común más corta entre s_ℓ y t_ℓ (resp. s_r y t_r). En un conjunto de contextos bi–direccionales cuyos elementos son no–disjuntos por pares también se define la extensión común más corta. El componente izquierdo (resp. derecho) es la extensión común más corta de los componentes izquierdos (resp. derechos) de los contextos del conjunto.

Ejemplo 4.16. La extensión común más corta de los contextos bi–direccionales $(01, 0)$ y $(0, 010)$ es $(01, 010)$.

Las relaciones prefijo y extensión son relaciones de orden parcial. Los bi–sufijos que determinan una clase \mathcal{C} de Π son sucesores (en la relación de orden extensión) de todos los contextos de \mathcal{C} . Por lo tanto el conjunto de sucesores de \mathcal{C} no es vacío. Con el siguiente lema establecemos que el conjunto de todos los sucesores de los elementos de la clase, tiene un mínimo, que es, por lo tanto, el supremo de la clase. Además, ese supremo pertenece a la clase, por lo que la clase tiene un contexto máximo.

Lema 4.17. Sea \mathcal{C} una clase de Π y sea (s_ℓ, s_r) el prefijo común más largo de $\mathcal{S}(\mathcal{C})$. Entonces, (s_ℓ, s_r) pertenece a \mathcal{C} y es la extensión común más corta de \mathcal{C} .

Demostración. Sea (t_ℓ, t_r) un contexto cualquiera de \mathcal{C} . Como s_ℓ es el prefijo común más largo de los componentes izquierdos de los bi–sufijos $\mathcal{S}(\mathcal{C})$ y t_ℓ es

también prefijo de esos componentes, se tiene que t_ℓ es prefijo de s_ℓ . De consideraciones análogas para los componentes derechos se llega a $\mathcal{C} \preceq (s_\ell, s_r)$, lo que prueba que (s_ℓ, s_r) es una extensión de todos los elementos de \mathcal{C} .

Entonces, por (4.2), $\mathcal{S}(s_\ell, s_r) \subseteq \mathcal{S}(\mathcal{C})$. Además se cumple que $(s_\ell, s_r) \preceq \mathcal{S}(\mathcal{C})$ y por lo tanto $\mathcal{S}(\mathcal{C}) \subseteq \mathcal{S}(s_\ell, s_r)$. De donde $\mathcal{S}(\mathcal{C}) = \mathcal{S}(s_\ell, s_r)$, lo que significa que (s_ℓ, s_r) pertenece a \mathcal{C} .

Queda por probar que (s_ℓ, s_r) es la más corta de las extensiones comunes de \mathcal{C} . Para esto basta observar que como (s_ℓ, s_r) pertenece a \mathcal{C} , la extensión común más corta de \mathcal{C} también debe ser extensión de (s_ℓ, s_r) . \square

El contexto (s_ℓ, s_r) definido en el Lema 4.17 será el elegido como el representativo de la clase y denotado por $R(\mathcal{C})$. Para cualquier contexto (t_ℓ, t_r) en \mathcal{C} , también se considera que (s_ℓ, s_r) es su representativo, lo cual se denota por $R(t_\ell, t_r)$. Para un contexto que no ocurre no se define su representativo. El conjunto de todos los representativos se denota por R .

Si $\mathcal{S}(\mathcal{C})$ tiene un único elemento, el contexto representativo de \mathcal{C} es ese único elemento. Hay n clases en estas condiciones, una para cada bi-sufijo. Estos contextos representativos se obtienen de manera trivial. En el caso uni-direccional ese único elemento de $\mathcal{S}(\mathcal{C})$ es una hoja del árbol de sufijos, y la clase la constituyen los contextos que son representados por la rama que entra en la hoja. Con el objetivo de obtener el contexto representativo de otras clases, presentaremos una propiedad que ese contexto debe cumplir.

Sea (s_ℓ, s_r) un contexto bi-direccional que no es bi-sufijo, y por lo tanto s_ℓ o s_r no son sufijos. Sin pérdida de generalidad supongamos que s_r no es sufijo. Entonces, para cada bi-sufijo en $\mathcal{S}(s_\ell, s_r)$ hay una única extensión mínima por derecha de (s_ℓ, s_r) que es prefijo de ese bi-sufijo. Por lo tanto, podemos particionar $\mathcal{S}(s_\ell, s_r)$ según las extensiones mínimas por derecha de (s_ℓ, s_r) que ocurren. Si esta partición tiene más de un elemento, decimos que (s_ℓ, s_r) es un *contexto ramificable por la derecha*. O sea que ocurren $(s_\ell, s_r a)$ y $(s_\ell, s_r b)$, con $a \neq b$. Es fácil ver que si esto sucede, s_r es un nodo interno del árbol de sufijos de y_r . De manera similar se define *contexto ramificable por la izquierda* y se ve que s_ℓ es nodo interno de y_ℓ .

Lema 4.18. *Si (s_ℓ, s_r) es ramificable por izquierda (resp. derecha) entonces s_r (resp. s_ℓ) no es sufijo.*

Demostración. Si (s_ℓ, s_r) es ramificable ocurre al menos dos veces. Pero si s_r fuese sufijo estaría en correspondencia uno a uno con una posición de x porque termina en $\#$. Por lo tanto no puede ser sufijo. \square

Como se mencionó en el caso uni–direccional, un nodo en un árbol de sufijos identifica unívocamente una clase \mathcal{C} de Π . Vamos a extender esa noción al caso bi–direccional.

Definición 4.19. Un contexto ramificable tanto por izquierda como por derecha se denomina *doblemente ramificable*, o *bi–ramificable*.

Un contexto bi–ramificable es prefijo de dos contextos bi–direccionales en los que es ramificable por izquierda y por derecha a la vez, como se ve en el siguiente lema.

Lema 4.20. *El contexto (s_ℓ, s_r) es bi–ramificable si y sólo si en la cadena de entrada x ocurren contextos $(s_\ell a, s_r b)$ y $(s_\ell c, s_r d)$ que cumplan $a, b, c, d \in \mathcal{A}$, $a \neq c$ y $b \neq d$.*

Demostración. La dirección *si* de la proposición se cumple por la Definición 4.19. Para demostrar la dirección *sólo si* veamos que por definición y por el Lema 4.18 como (s_ℓ, s_r) es ramificable por izquierda ocurren los bi–contextos $(s_\ell a, s_r b)$ y $(s_\ell a', s_r b')$, con $a \neq a'$. Si, además $b \neq b'$, la proposición queda probada. Por lo tanto suponemos $b = b'$.

Como (s_ℓ, s_r) es también ramificable por derecha, ocurren los bi–contextos $(s_\ell c', s_r d')$ y $(s_\ell c, s_r d)$, con $d \neq d'$. Por la misma razón que antes, suponemos $c = c'$.

Como $a \neq a'$, c es diferente de al menos uno de ellos. Sin pérdida de generalidad, suponemos $a \neq c$. Por la misma razón, b debe ser diferente a al menos uno de $\{d, d'\}$, y sin pérdida de generalidad elegimos $b \neq d$. De esta manera queda probada la proposición. \square

El conjunto de extensiones bi–ramificables de (s_ℓ, s_r) se denota por $\mathcal{B}(s_\ell, s_r)$. Y \mathcal{B} denota a todos los bi–ramificables, o sea $\mathcal{B}(\epsilon, \epsilon)$. Por lo que hemos visto, \mathcal{B} es un subconjunto del producto cartesiano de los nodos internos de los árboles de sufijos correspondientes a y_ℓ y y_r .

Una extensión mínima $(s_\ell a, s_r)$, que ocurre, de un contexto bi–ramificable (s_ℓ, s_r) no puede pertenecer a la misma clase que (s_ℓ, s_r) , porque el conjunto $\mathcal{S}(s_\ell a, s_r)$ está contenido de manera estricta en $\mathcal{S}(s_\ell, s_r)$. Por lo tanto los contextos de $\mathcal{C}(s_\ell, s_r)$ deben ser prefijos de (s_ℓ, s_r) . En el Lema 4.17 habíamos establecido que lo que después denominamos representativo es la extensión de todos los contextos de la clase. Esto sugiere que los contextos bi–ramificables son representativos. Esto quedará demostrado con los Lemas 4.23 y 4.24.

Ejemplo 4.21. Siguiendo con el Ejemplo 4.3, $(01, 010)$ es el representativo de su clase y podemos ver que ocurren sus extensiones $(010, 0101)$ en la posición 7 y $(011, 0100)$ en la posición 5.

El Lema 4.23 relaciona contextos representativos con bi-sufijos y contextos bi-ramificables: $R \subseteq \mathcal{S} \cup \mathcal{B}$.

Antes vemos un lema previo.

Lema 4.22. *El prefijo común más largo de un conjunto de contextos uni-direccionales es el prefijo común más largo de al menos un par de los contextos del conjunto.*

Demostración. Alcanza con mostrar un ejemplo: el par formado por el primero y el último en orden lexicográfico de los contextos del conjunto. \square

Lema 4.23. *Un contexto representativo es o bien bi-sufijo, o bien bi-ramificable.*

Demostración. Un contexto representativo es la extensión común más corta de los elementos de una clase \mathcal{C} de Π y el prefijo común más largo de $\mathcal{S}(\mathcal{C})$. Sea (s_ℓ, s_r) el contexto representativo de \mathcal{C} .

Si $\mathcal{S}(\mathcal{C})$ tiene un solo bi-sufijo, entonces (s_ℓ, s_r) es ese bi-sufijo por ser trivialmente el prefijo común más largo. Si, en cambio, más de un bi-sufijo pertenece a $\mathcal{S}(\mathcal{C})$, entonces, por el Lema 4.22 hay al menos un par de bi-sufijos de $\mathcal{S}(\mathcal{C})$ para los cuales el prefijo común más largo de sus componentes izquierdos es s_ℓ . Entonces, como (s_ℓ, s_r) es prefijo de ambos bi-sufijos, es ramificable por izquierda. De manera similar se ve que es ramificable por derecha, por lo que es bi-ramificable. \square

El Lema 4.23 establece una condición necesaria para que un contexto sea representativo. Del Lema 4.24 se deduce que sólo un contexto en cada clase satisface esta condición. Por lo tanto la condición es suficiente.

Lema 4.24. *Sean (s_ℓ, s_r) y (t_ℓ, t_r) dos contextos diferentes, cada uno de ellos bi-sufijo o bi-ramificable. Entonces, $\mathcal{S}(s_\ell, s_r) \neq \mathcal{S}(t_\ell, t_r)$.*

Demostración. Si las extensiones de los dos contextos son conjuntos disjuntos la proposición se cumple. En otro caso, o bien t_ℓ es prefijo de s_ℓ , o bien s_ℓ es prefijo de t_ℓ . Algo similar ocurre con los contextos derechos. Como los dos contextos son diferentes, en alguno de los dos casos la relación es estricta. Sin pérdida de generalidad supongamos $t_\ell \prec s_\ell$. Esto significa que $t_\ell a \preceq s_\ell$ para algún $a \in \mathcal{A}$. Además (t_ℓ, t_r) no puede ser bi-sufijo y por lo tanto es bi-ramificable. Esto significa que debe ocurrir $(t_\ell b, t_r)$, con $b \neq a$. Por lo tanto $\mathcal{S}(t_\ell b, t_r)$ no es vacío, y es disjunto con $\mathcal{S}(t_\ell a, s_r)$, de donde $\mathcal{S}(t_\ell b, t_r) \cap \mathcal{S}(s_\ell, s_r) = \emptyset$, porque $\mathcal{S}(s_\ell, s_r) \subseteq \mathcal{S}(t_\ell a, s_r)$. Y como $\mathcal{S}(t_\ell b, t_r) \subset \mathcal{S}(t_\ell, t_r)$ se concluye $\mathcal{S}(s_\ell, s_r) \neq \mathcal{S}(t_\ell, t_r)$. \square

Los Lemas 4.17, 4.23, y 4.24 implican el siguiente teorema, que establece $R = \mathcal{S} \cup \mathcal{B}$.

Teorema 4.25. *El conjunto de contextos representativos es la unión del conjunto de bi-sufijos y el conjunto de contextos bi-ramificables.*

4.3.1. Conjunto óptimo

Ahora aplicamos la noción de contexto representativo para la búsqueda de un conjunto óptimo de contextos para una función de peso dada. En la Sección 3.2, habíamos dado la Definición 3.8 de conjunto válido relativo a un contexto. Vamos a adaptarla para tratar sólo con contextos representativos. Es evidente que uno de los conjuntos válidos de contextos relativo a (s_ℓ, s_r) es el conjunto de sus extensiones mínimas izquierdas, $\{s_\ell a, s_r : a \in \mathcal{A}\}$. Algunas de esas extensiones no ocurren y otras no son contextos representativos. Ahora vemos otro conjunto válido de gran utilidad. Está compuesto por los $(s_\ell a, s_r)$ que no ocurren y los representativos de los $(s_\ell a, s_r)$ que ocurren. Es decir, es la unión de

$$\{(s_\ell a, s_r) : a \in \mathcal{A}, \mathcal{I}(s_\ell a, s_r) = \emptyset\} \text{ y } \{R(s_\ell a, s_r) : a \in \mathcal{A}, \mathcal{I}(s_\ell a, s_r) \neq \emptyset\}.$$

El primero de los conjuntos no ocurre, por lo que no hay bi-sufijos que sean extensión suya. Si $(s_\ell a, s_r)$ ocurre, entonces $\mathcal{S}(s_\ell a, s_r)$ no es vacío y es igual a $\mathcal{S}(R(s_\ell a, s_r))$. Y cada contexto de $\mathcal{S}(s_\ell, s_r)$ es extensión de exactamente un elemento del segundo conjunto. Esto motiva a dar la siguiente definición.

Definición 4.26. Un conjunto de contextos \mathcal{V} es *compacto relativo* a (s_ℓ, s_r) , si es extensión de (s_ℓ, s_r) , cada uno de sus elementos es representativo, y cada elemento de $\mathcal{S}(s_\ell, s_r)$ es extensión de exactamente un elemento de \mathcal{V} .

Debe quedar claro que en esta definición el contexto (s_ℓ, s_r) no tiene por qué ser representativo.

Definición 4.27. Un conjunto de contextos \mathcal{V} es *óptimo compacto relativo* a (s_ℓ, s_r) si tiene peso mínimo entre todos los conjuntos compactos relativos a (s_ℓ, s_r) , donde el peso de un conjunto de contextos es definido como la suma de los pesos de sus elementos.

Nota 4.28. Se debe notar que dado un conjunto óptimo compacto relativo a (s_ℓ, s_r) se puede construir un conjunto que lo contenga y sea óptimo relativo a (s_ℓ, s_r) (según la Definición 3.8 y siguientes). Cada contexto que se agrega no ocurre y por lo tanto su vector de conteos es nulo así como también lo es su peso.

Teorema 4.29. *Sea \mathcal{C} una clase de Π . Entonces, un conjunto óptimo compacto de contextos relativo a $R(\mathcal{C})$ es también óptimo compacto relativo a cualquier otro contexto en \mathcal{C} .*

Demostración. Alcanza con ver que cualquier conjunto compacto relativo a cualquier (s_ℓ, s_r) perteneciente a \mathcal{C} es también compacto relativo a $R(\mathcal{C})$. Y esto se cumple porque cualquier contexto representativo que sea extensión de (s_ℓ, s_r) también lo es de $R(\mathcal{C})$. \square

De este teorema y de la Nota 4.28 se deriva el siguiente corolario.

Corolario 4.30. *El problema de encontrar un conjunto óptimo de contextos es equivalente a encontrar un conjunto óptimo compacto.*

Debe notarse que un conjunto óptimo compacto es, según la Definición 4.26 un subconjunto de los contextos representativos, que a su vez, según el Teorema 4.25 son los contextos bi-sufijos y bi-ramificables.

Bases para el nuevo algoritmo Ahora vamos a reformular expresiones basándonos en el actual marco de sólo trabajar con contextos representativos.

De la Ecuación (3.3) se sabe que el peso óptimo de un contexto (s_ℓ, s_r) se obtiene del mínimo entre tres valores: la función $w(s_\ell, s_r)$, la suma de los pesos óptimos de las extensiones mínimas izquierdas y la suma de los pesos óptimos de las extensiones mínimas derechas. Para cada contexto representativo (s_ℓ, s_r) , un conjunto válido relativo a él, trivial, es el propio (s_ℓ, s_r) . Ahora definimos otros dos conjuntos válidos relativos a (s_ℓ, s_r) .

Definición 4.31. *El conjunto válido mínimo por izquierda (resp. derecha) relativo a (s_ℓ, s_r) consiste en los contextos representativos de las extensiones mínimas izquierdas (resp. derechas) que ocurren.*

O sea, son los conjuntos

$$\{R(s_\ell a, s_r) : a \in \mathcal{A}, \mathcal{I}(s_\ell a, s_r) \neq \emptyset\} \quad (4.3)$$

y

$$\{R(s_\ell, s_r a) : a \in \mathcal{A}, \mathcal{I}(s_\ell, s_r a) \neq \emptyset\}. \quad (4.4)$$

Primero vamos a ver que el vector de conteos de (s_ℓ, s_r) se puede calcular sumando los vectores de conteos de cualquiera de estos conjuntos. La Ecuación (3.4) establece que el vector de conteos es la suma de los vectores de conteos de las extensiones mínimas izquierdas o derechas.

Esa suma se puede separar en dos, una con las extensiones que no ocurren y otra con las que ocurren:

$$c(s_\ell, s_r) = \sum_{a \in \mathcal{A}} c(s_\ell a, s_r) = \sum_{\substack{a \in \mathcal{A}, \\ \mathcal{I}(s_\ell a, s_r) = \emptyset}} c(s_\ell a, s_r) + \sum_{\substack{a \in \mathcal{A}, \\ \mathcal{I}(s_\ell a, s_r) \neq \emptyset}} c(s_\ell a, s_r).$$

Como los vectores de los contextos que no ocurren son nulos, la primera suma es 0. El vector de conteos de un contexto es igual al de su representativo porque ocurren en las mismas posiciones. Por lo tanto en la segunda suma los contextos pueden reemplazarse por sus representativos. Entonces (y procediendo de manera similar con las extensiones mínimas por derecha) se llega a una nueva expresión para calcular el vector de conteos de un contexto si se conoce los vectores de conteos de sus extensiones.

$$\mathbf{c}(\mathbf{s}_\ell, \mathbf{s}_r) = \sum_{\substack{a \in \mathcal{A}, \\ \mathcal{I}(\mathbf{s}_\ell a, \mathbf{s}_r) \neq \emptyset}} \mathbf{c}(\mathbf{R}(\mathbf{s}_\ell a, \mathbf{s}_r)) = \sum_{\substack{a \in \mathcal{A}, \\ \mathcal{I}(\mathbf{s}_\ell, \mathbf{s}_r a) \neq \emptyset}} \mathbf{c}(\mathbf{R}(\mathbf{s}_\ell, \mathbf{s}_r a)). \quad (4.5)$$

Finalmente adaptamos la optimalidad expresada en (3.3). El razonamiento es similar al anterior. Cada una de las sumas se separa en dos. Una incluye las extensiones mínimas que no ocurren, que tienen peso óptimo 0, y por lo tanto esa suma es 0. La otra suma incluye las extensiones que ocurren. Del Teorema 4.29 deducimos que el peso óptimo de un contexto es igual al de su representativo. Por lo tanto hacemos la sustitución de cada extensión mínima con su representativo. Por otro lado, el cálculo de $\mathbf{w}(\mathbf{s}_\ell, \mathbf{s}_r)$ toma como parámetro el vector de conteos, que vimos que se calcula correctamente, por lo que no es afectado por la restricción de sólo tratar con contextos representativos. Entonces, definiendo

$$\begin{aligned} N &= \mathbf{w}(\mathbf{s}_\ell, \mathbf{s}_r) \\ L &= \sum_{\substack{a \in \mathcal{A}, \\ \mathcal{I}(\mathbf{s}_\ell a, \mathbf{s}_r) \neq \emptyset}} \mathbf{w}_{opt}(\mathbf{R}(\mathbf{s}_\ell a, \mathbf{s}_r)) \\ R &= \sum_{\substack{a \in \mathcal{A}, \\ \mathcal{I}(\mathbf{s}_\ell, \mathbf{s}_r a) \neq \emptyset}} \mathbf{w}_{opt}(\mathbf{R}(\mathbf{s}_\ell, \mathbf{s}_r a)) \end{aligned}$$

llegamos a la expresión de optimalidad entre contextos representativos:

$$\mathbf{w}_{opt}(\mathbf{s}_\ell, \mathbf{s}_r) = \min \{N, L, R\}. \quad (4.6)$$

Y entonces se puede decidir cuál es la partición que debe aplicarse a $(\mathbf{s}_\ell, \mathbf{s}_r)$ para obtener un conjunto óptimo relativo a él.

$$\mathbf{split}(\mathbf{s}_\ell, \mathbf{s}_r) = \mathbf{argmin} \{N, L, R\} \quad (4.7)$$

donde la función **argmin** devuelve ‘N’, ‘L’ o ‘R’ si el mínimo de $\{N, L, R\}$ es N , L o R respectivamente y resolviendo los empates de alguna manera determinística.

Capítulo 5

Subárboles Compactos

En este capítulo vamos a describir un método eficiente para hallar todos los contextos representativos correspondientes a una secuencia x . El método utiliza los árboles de sufijos de las cadenas formadoras de contextos. Comienza obteniendo los bi-sufijos y luego recorre uno de los árboles en post-orden, obteniendo para cada nodo todos los nodos del otro árbol con los cuales forma un contexto representativo. Se verá que la cantidad de contextos representativos es $O(n \cdot h)$, siendo h el mínimo de las alturas de los dos árboles y n el largo de la secuencia x .

Por cada nodo del árbol recorrido se encontrará un conjunto de nodos del otro árbol, al que se le puede dar una estructura de árbol. Se van a definir operaciones entre estos árboles. En el Capítulo 6 se define una clase de grafos y se describe su construcción en la Sección 6.2. Esa construcción puede verse como una generalización de las operaciones sobre árboles que aquí se definen.

Vamos a dar una definición de lca , el ancestro común más profundo de los nodos de un árbol, que es una importante operación entre nodos de un árbol. En el Apéndice A veremos que se puede calcular en tiempo constante.

5.1. Árboles

En esta sección trataremos los árboles de manera estructural, sin imponer ningún tipo de etiquetas en las ramas. Vamos a definir dos operaciones entre árboles que en la Sección 5.2 serán usadas para encontrar los contextos representativos.

Sea T un árbol. Denotamos con $\mathbf{V}(T)$ al conjunto de los nodos del árbol, que es la unión disjunta del conjunto de sus nodos internos, denotado $\mathbf{B}(T)$, y el conjunto de sus hojas, denotado $\mathbf{L}(T)$. Si T es el árbol vacío, cada uno de estos conjuntos es vacío. Si el nodo v es ancestro del nodo w lo denotamos por $v \preceq w$, y en otro caso por $v \not\preceq w$. Si w es hijo de v la rama que conecta v y w se denota

\overrightarrow{vw} , tal como se hizo al definir árboles en el Capítulo 2, aunque en esta sección las ramas no están etiquetadas. La función $\text{root}(T)$ identifica el nodo que es la raíz del árbol T . Y $\text{children}(T, v)$, con $v \in \mathbf{V}(T)$, es el conjunto de nodos de T que son los hijos del nodo v . El árbol vacío se representa con \perp .

Un *árbol compacto* es un árbol en que cada nodo interno tiene dos o más hijos. En esta tesis sólo se consideran árboles compactos.

5.1.1. Subárboles Compactos

Definición 5.1. Sean v y w nodos de un árbol B , con $v \preceq w$. Se dice que \overrightarrow{vw} es la *compresión del camino* $\overrightarrow{u_1 u_2} \overrightarrow{u_2 u_3} \cdots \overrightarrow{u_{n-1} u_n}$ de B , donde $v = u_1, w = u_n$, y $\overrightarrow{u_i u_{i+1}}$ es una rama de B , con $1 \leq i \leq n - 1$.

Definición 5.2. Los caminos \overrightarrow{vy} y \overrightarrow{vz} de un árbol son disjuntos si no tienen ramas en común.

Esto significa que si \overrightarrow{vy} y \overrightarrow{vz} son disjuntos, entonces $\overrightarrow{vy} = \overrightarrow{vu_2} \cdots \overrightarrow{u_n y}$ y $\overrightarrow{vz} = \overrightarrow{vw_2} \cdots \overrightarrow{w_k z}$, con $u_2 \neq w_2$.

Definición 5.3. Un *subárbol compacto*, T , del árbol compacto B , al que denominamos *árbol base*, es un árbol compacto cuyos nodos pertenecen a $\mathbf{V}(B)$, las ramas son caminos de B , y todo par de ramas salientes de un nodo de T son disjuntas al ser consideradas como caminos de B .

Un ejemplo se ve en la Figura 5.1.

Según la definición, las ramas del subárbol no tienen por qué ser ramas del árbol base. Como ejemplo, en la Figura 5.1, la rama $\overrightarrow{\gamma\delta}$ en T es la compresión del camino $\overrightarrow{\gamma\delta} \overrightarrow{\delta\epsilon}$ en B .

De la definición de subárbol compacto se deduce que los hijos de un nodo u en el subárbol son descendientes de hijos diferentes del mismo nodo u en el árbol base.

Entonces, formalmente T es un subárbol compacto de B si y sólo si

$$\mathbf{V}(T) \subseteq \mathbf{V}(B),$$

Si $v \in \mathbf{B}(T)$ entonces $|\text{children}(T, v)| \geq 2$,

Si $v, w \in \text{children}(T, u), v \neq w$ entonces

$$\exists v', w' \in \text{children}(B, u), v' \neq w', \text{ tal que } v' \preceq v, w' \preceq w.$$

Denotamos con B^C a la colección de subárboles compactos del árbol base B .

Si u es un nodo de B , con $B(u)$ denotamos el subárbol compacto de B formado por todos los nodos de B que descienden de u . Obviamente, la raíz de este subárbol es u .

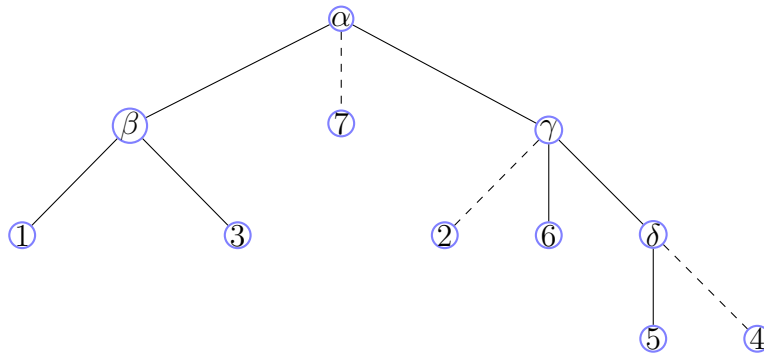


Figura 5.1: Subárbol compacto. En trazo continuo se ve un subárbol compacto sobre un árbol base que se completa en trazo discontinuo. Las hojas 7, 2 y 4 y el nodo interno δ del árbol base no son nodos del subárbol. Hay que notar que si δ perteneciera al subárbol, este no sería compacto. La rama del subárbol que va desde γ hasta 5 es la compresión del camino que va desde γ hasta δ y desde δ hasta 5 en el árbol base.

Sea T un subárbol compacto de B y v un nodo de B que es descendiente de la raíz de T . Más formalmente, $T \in B^C$, $v \in \mathbf{V}(B)$, $\text{root}(T) \preceq v$. La definición de $T(v)$ es igual a la de $B(v)$, es decir, es el subárbol formado por los nodos de T que descienden de v . Es posible que v no pertenezca a T , y esto genera dos casos. En uno de ellos ningún nodo de T desciende de v , por lo que $T(v)$ es \perp . En el otro caso todos los nodos de T que descienden de v lo hacen por la misma rama de T , siendo la raíz de $T(v)$ el menos profundo de esos nodos.

Ejemplo 5.4. En la Figura 5.1 $T(\beta)$ es igual a $B(\beta)$. $T(\gamma)$ es el subárbol con los nodos γ , 6 y 5. Los dos casos en los que v no pertenece a T se ejemplifican con $T(\delta)$, que es solamente la hoja 5 y $T(4)$ que es \perp .

5.1.2. LCA

El lca de dos nodos de un árbol es el ancestro común de ambos que está a mayor profundidad. De manera obvia extendemos la definición de lca a un conjunto de nodos.

Definición 5.5. El lca (“lowest common ancestor”) de un conjunto de dos o más nodos de un árbol T es el nodo que está a mayor profundidad en el árbol entre los que son ancestros a todos los nodos del conjunto.

Formalmente,

$$\text{lca} : \text{set}(\mathbf{V}(T)) \rightarrow \mathbf{B}(T)$$

$$\text{Precondición} : |S| \geq 2,$$

$$\text{lca}(S) \preceq S,$$

$$v \in \mathbf{V}(T), v \preceq S \Rightarrow v \preceq \text{lca}(S).$$

Ejemplo 5.6. En la Figura 5.1 β es el lca de los nodos 1 y 3; α es el lca de los nodos 1, 7 y γ ; γ es el lca de 2, γ y 6.

Nota 5.7. El lca de dos elementos del conjunto S no necesariamente es el lca del conjunto. Lo que sí se cumple es que hay al menos un par de elementos de S cuyo lca es el lca del conjunto. Esto es trivial si $\text{lca}(S)$ pertenece a S . En el otro caso, si no fuera así significaría que todos los nodos del conjunto descienden de $\text{lca}(S)$ por la misma rama, por lo que el hijo de $\text{lca}(S)$ según esa rama sería un ancestro común más profundo de los nodos de S , en contra de lo supuesto. Un ejemplo es el par formado por el primero y el último de los nodos de S en una recorrida euleriana (“depth first”). Esto tiene importancia desde el punto de vista de la eficiencia porque seleccionando adecuadamente un par de elementos se puede obtener el lca de todo el conjunto. Aunque en esta sección los nodos de los árboles no están asociados a contextos se debe notar la relación con el Lema 4.22.

Ejemplo 5.8. El lca de los nodos 1, 3 y γ es α que también es el lca de 1 y γ (aunque no es el lca de los nodos 1 y 3); el lca de 2, γ y 6 es γ que es el lca de 2 y 6.

5.1.3. LCA generalizado

El lca es un nodo. Según lo que se estableció en la Nota 5.7 puede haber más de un nodo que sea lca de algún par de nodos del conjunto S . Por lo tanto podemos definir una nueva operación que generalice lca y devuelva el conjunto de nodos que son lca de algún par de esos nodos. En el caso particular en que el lca de cada par de nodos de S sea el mismo la nueva operación es la operación lca ya definida.

Vamos a extender esa generalización en otro aspecto. El lca se define sobre el conjunto de nodos. Pero como un nodo es la raíz de un subárbol, la definición podría haberse hecho sobre el conjunto de subárboles. Al considerarlo de esta manera la generalización de la operación puede obtener no sólo él o los lca de las raíces, sino el de todos los posibles pares de nodos tomados de diferentes subárboles.

Ejemplo 5.9. En el árbol base de la Figura 5.1 consideremos los subárboles definidos por los nodos $\{\alpha, 3, 6\}$ y $\{\delta, 4, 5\}$. La generalización del lca es el conjunto $\{\alpha, \gamma\}$. Esto es porque el lca de los nodos α o 3 con cualquiera de los del otro conjunto es α , mientras que el lca del nodo 6 con los otros es γ .

Aquí generalizamos el concepto lca para considerar una colección de conjuntos de nodos.

Definición 5.10. El *lca generalizado*, glca , de una colección de conjuntos de nodos de un árbol es el conjunto de nodos formado por el lca de todos los pares de nodos tomados de distintos conjuntos. Si hay menos de dos conjuntos de nodos, el resultado es \emptyset .

Nota 5.11. Un mismo nodo v puede pertenecer a varios conjuntos de la colección. En ese caso, $v = \text{lca}(v, v)$ pertenece al glca .

Lo que nos interesa es la restricción de esta operación cuando cada conjunto de nodos tiene estructura de subárbol compacto. La importancia de restringir el glca a subárboles compactos proviene de que, como veremos en la Sección 5.2.1, con esta operación se obtienen los contextos representativos. Sólo trataremos con subárboles compactos cuyas hojas son hojas del árbol base, y además, los conjuntos de hojas son disjuntos por pares.

La siguiente es la especificación de glca cumpliendo estas restricciones:

$$\text{glca} : \text{set}(B^C) \rightarrow \text{set}(\mathbf{B}(B))$$

Precondiciones :

$$\mathbf{L}(T) \subseteq \mathbf{L}(B), \text{ para todo } T \in S,$$

$$\mathbf{L}(T_1) \cap \mathbf{L}(T_2) = \emptyset, \text{ con } T_1, T_2 \in S, T_1 \neq T_2,$$

$$\text{glca}(S) = \begin{cases} \emptyset & \text{si } |\{T \neq \perp : T \in S\}| < 2 \\ V & \text{en otro caso, donde} \\ & v \in V \Leftrightarrow \exists v_i \in \mathbf{V}(T_i), v_j \in \mathbf{V}(T_j) \text{ tal que } v = \text{lca}(v_i, v_j), \\ & \text{con } T_i, T_j \in S, T_i \neq T_j. \end{cases}$$

Es posible obtener el glca de manera eficiente. Primero se puede ver que no es necesario obtener el lca de cada par de nodos. Se obtiene el mismo resultado restringiendo la operación a los pares de hojas, ya que la última expresión puede sustituirse por

$$v \in V \Leftrightarrow \exists l_i \in \mathbf{L}(T_i), l_j \in \mathbf{L}(T_j) \text{ tal que } v = \text{lca}(l_i, l_j), \\ \text{con } T_i, T_j \in S, T_i \neq T_j,$$

porque para $v_i \in \mathbf{V}(T_i)$ existe $l_i \in \mathbf{L}(T_i)$ que cumple $v_i \preceq l_i$ y para $v_j \in \mathbf{V}(T_j)$ existe $l_j \in \mathbf{L}(T_j)$ que cumple $v_j \preceq l_j$, y $\mathbf{lca}(v_i, v_j) = \mathbf{lca}(l_i, l_j)$.

Aún así, ese procedimiento es ineficiente porque algunos nodos del \mathbf{glca} son el \mathbf{lca} de varios pares de hojas y se obtendrían varias veces.

Vamos a demostrar que la siguiente especificación, que elimina las ineficiencias mencionadas, describe el \mathbf{glca} .

$$\mathbf{glca}(S) = \begin{cases} \emptyset & \text{si } |S| < 2 \\ \{v\} \cup_{w \in \mathbf{children}(B, v)} \mathbf{glca} \{T(w) : T \in S, T(w) \neq \perp\} & \\ \text{en otro caso, donde } v = \mathbf{lca} \{\mathbf{root}(T) : T \in S\}. & \end{cases} \quad (5.1)$$

Lema 5.12. *La Expresión (5.1) es el \mathbf{glca} del conjunto S de subárboles compactos del árbol base B .*

Demostración. Está claro que todo nodo obtenido usando esta especificación pertenece al \mathbf{glca} porque es el resultado de aplicar el \mathbf{lca} a nodos de diferentes subárboles. Para ver el recíproco, esto es, que el \mathbf{lca} de todo par de hojas de diferentes subárboles pertenece al conjunto obtenido, sea $z = \mathbf{lca}(l_i, l_j)$, donde l_i y l_j son respectivamente hojas de los subárboles T_i y T_j de S , con $T_i \neq T_j$. Hay dos casos posibles. Uno es que z sea igual a v . Si es así el \mathbf{lca} de las dos hojas pertenece al resultado obtenido. En el otro caso, z pertenece a $B(w)$ para algún w en $\mathbf{children}(B, v)$, y $T_i(w)$ y $T_j(w)$ son no vacíos. Entonces, l_i y l_j son hojas de diferentes subárboles en la obtención del \mathbf{glca} de los subárboles que descienden de w . \square

Esta especificación induce de manera directa un algoritmo recursivo para obtener el \mathbf{glca} . Como se ve, cada nodo del \mathbf{glca} se obtiene sólo una vez, por lo que, la complejidad del algoritmo, medida en cantidad de veces que se invoca la función \mathbf{lca} es igual a la cantidad de nodos del \mathbf{glca} . En el Apéndice A se verá que el cálculo del \mathbf{lca} se puede hacer en $O(1)$. Como consecuencia, la complejidad de obtener el \mathbf{glca} es $O(|\mathbf{glca}|)$.

Para consideraciones de búsqueda de eficiencia en los algoritmos notemos algunos casos particulares de $T(w)$ en la Ecuación (5.1).

$$T(w) = \begin{cases} T & \text{si } w \preceq \mathbf{root}(T), \\ \perp & \text{si } w' \preceq \mathbf{root}(T) \\ & \text{para algún } w' \in \mathbf{children}(B, v), w' \neq w, \\ \neq T, \text{ tal vez } \perp & \text{si } v = \mathbf{root}(T). \end{cases}$$

En el tercer caso la raíz del árbol T es el glca de S . En ese caso es incierto en qué contribuye T a la llamada recursiva. El aporte es un subárbol estricto de T , tal vez \perp . Si, en cambio, $v \neq \text{root}(T)$ estamos en los otros dos casos. En el primero, el árbol completo es descendiente de w por lo que todos sus nodos quedan involucrados en los cálculos. En el segundo caso ningún nodo de T desciende de w .

5.1.4. Árbol soporte

Un subárbol compacto del árbol base B queda determinado por sus hojas. Esto es porque cualquier nodo interno es el lca de algún par de hojas que descienden de él según ramas diferentes. Recíprocamente, el lca de cualquier par de hojas es un nodo interno. El subárbol es la clausura del conjunto de hojas bajo la operación lca . Esto es, empezando con un conjunto que contiene las hojas, se obtiene el lca de un par de elementos del conjunto y el resultado se agrega al conjunto. El proceso termina cuando todo posible par da como resultado un elemento que ya está en el conjunto. Entonces, es el mínimo subárbol compacto que contiene todo el conjunto de hojas. No cualquier conjunto de nodos de B puede ser un conjunto de hojas de un subárbol compacto. Nos limitaremos a subárboles cuyas hojas sean hojas de B . Con esta restricción vamos a dar la definición de árbol soporte.

Definición 5.13. El *árbol soporte* de un conjunto de hojas L del árbol base B , es el subárbol compacto de B cuyas hojas son exactamente L .

Podemos interpretar el árbol soporte como una función sobre conjuntos de hojas:

$$\begin{aligned} \text{soporte} : \text{set}(\mathbf{L}(B)) &\rightarrow B^C \\ \text{soporte}(L) &= T \text{ tal que} \\ \mathbf{V}(T) &= \bigcup_{l_i, l_j \in L} \text{lca}(l_i, l_j) \end{aligned}$$

Se cumple que $\mathbf{L}(T)$ es L . Estas hojas se obtienen cuando $l_i = l_j$. Si L es vacío el resultado es el árbol vacío.

Ejemplo 5.14. En la Figura 5.1 lo que se ve en trazo continuo es el árbol soporte de las hojas 1, 3, 6 y 5.

5.1.5. Unión de subárboles

Para el proceso de obtención de los contextos representativos vamos a definir otra operación entre subárboles compactos. Esta definición está restringida a

subárboles para los cuales el conjunto de hojas de todos ellos satisface la propiedad de que dado cualquier par de esas hojas de los subárboles, al ser consideradas como nodos del árbol base, ninguna es ancestro de la otra. Esto es, para todo par v_1, v_2 de hojas de los subárboles se cumple que en B , $v_1 \not\prec v_2$, $v_2 \not\prec v_1$. Esto se cumple con las restricciones que hemos impuesto, esto es, que las hojas de los subárboles son hojas del árbol base y los conjuntos de hojas de cada par de subárboles son disjuntos.

Definición 5.15. Llamamos *unión de subárboles compactos* de un árbol base al árbol soporte de la unión de los conjuntos de hojas de los subárboles.

Formalmente:

$$\text{union} : \text{set}(B^C) \rightarrow B^C$$

$$\text{Precondición} : |S| \geq 2$$

$$\text{union}(S) = \text{soporte} \left(\bigcup_{T \in S} L(T) \right).$$

Con la condición impuesta de que las hojas de los subárboles son hojas del árbol base, esta operación es cerrada, es decir, el resultado es también un subárbol cuyas hojas son hojas del árbol base.

Como ya fue establecido al definir árbol soporte, es posible obtener la unión mediante la clausura lca del conjunto de hojas. Además habrá que establecer las ramas. Una primera simplificación de este algoritmo es notar que, con el lca de dos hojas del mismo subárbol se obtiene un nodo interno que ya se conocía. Entonces, para obtener todos los nodos de la unión basta con restringir el cálculo a hojas de diferentes subárboles. Y esto es obtener el glca de los subárboles. O sea que

$$\mathbf{V}(\text{union}(S)) = \text{glca}(S) \cup \bigcup_{T \in S} \mathbf{V}(T). \quad (5.2)$$

Con el procedimiento descrito es posible obtener repetidas veces el mismo nodo. Una versión más eficiente es una variante de la solución vista para el glca en la Ecuación (5.1). Se describe en el Algoritmo 2. La generalización para el caso bi-direccional, en el que se obtiene un grafo, se verá en el Algoritmo 5

Se va a demostrar la correctitud de este algoritmo.

Teorema 5.16. *El Algoritmo 2 devuelve la unión de los subárboles pertenecientes al conjunto S .*

Demostración. Por la Ecuación (5.2) hay que demostrar que la estructura devuelta, U , contiene todos los nodos del glca y de todos los árboles de S . Además hay que probar que la estructura de ramas es correcta.

Algoritmo 2 $\text{union}(S : \text{set}(B^C))$

```
1: Salida:  $U : B^C$ 
2: Precondición:  $|S| \geq 2$ 
3:  $v = \text{root}(U) \leftarrow \text{lca} \{ \text{root}(T) : T \in S \}$ 
4: Para cada  $w \in \text{children}(B, v)$ 
5:    $S' \leftarrow \emptyset$ 
6:   Para cada  $T \in S$ 
7:     Si  $T(w) \neq \perp$  entonces
8:        $S' \leftarrow S' \cup \{T(w)\}$ 
9:   Si  $S' \neq \emptyset$  entonces
10:     $U(w) \leftarrow \begin{cases} T' & \text{si } S' = \{T'\} \\ \text{union}(S') & \text{si } |S'| \geq 2. \end{cases}$ 
11: devuelve  $U$ 
```

Como puede verse comparando con la Ecuación (5.1) los nodos del glca son los que resultan de obtener en la línea 3 el lca de las raíces de los subárboles en todas las llamadas recursivas de la línea 10.

Para cada $T \in S$, su raíz, $\text{root}(T)$, pertenece a U si es igual a v . En otro caso $\text{root}(T)$ desciende de algún hijo w de v por lo que es raíz de $T(w)$. Como tal, en la línea 10 o bien, si es el único árbol en S' , es incluido en U , junto a todos sus descendientes, o bien es la raíz de uno de los árboles de la llamada recursiva, en una instancia con árboles más chicos, por lo que en alguna de las llamadas va a ser incluido en U . Volviendo al primer caso, $\text{root}(T) = v$, cada hijo de $\text{root}(T)$ es raíz de alguno de los $T(w)$, por lo que valen para ellos las consideraciones anteriores y serán incluidos en U .

En cuanto a la estructura de ramas, basta ver que $U(w)$, el subárbol de U , que desciende de v por la misma rama que une a v con w está formado sólo por descendientes de w , ya que por definición esto lo cumplen todos los árboles $T(w)$. \square

La complejidad de este algoritmo, medida en la cantidad de veces que se invoca la función lca , es la misma que la de obtener el glca . Para ello se impone como precondición que el conjunto S tenga más de un elemento para evitar las recorridas de los árboles que ya fueron construidos.

Usaremos el árbol T_{y_ℓ} de la Figura 4.2 como ejemplo de árbol base. En la Figura 5.2, en la parte de arriba, vemos tres subárboles de T_{y_ℓ} . Con trazo discontinuo aparece el árbol base. Dos de los subárboles son hojas de T_{y_ℓ} . La raíz de la unión es el lca de las raíces de los tres subárboles, que es la raíz de T_{y_ℓ} . Descendiendo de esta raíz por la rama izquierda sólo está el tercer subárbol, por lo cual se establece una rama hacia él. Nada desciende por la rama del medio. Por la

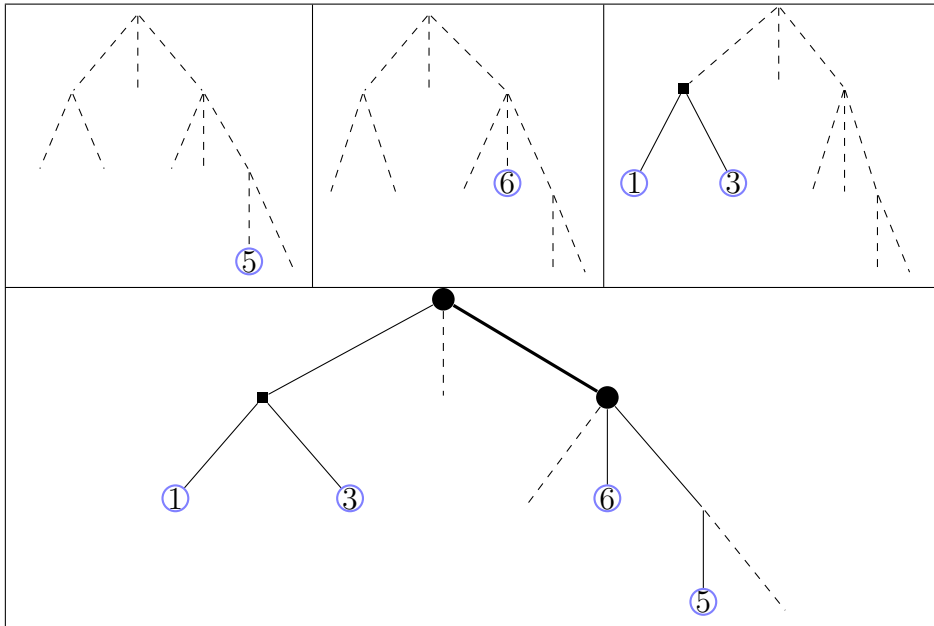


Figura 5.2: Unión de subárboles

rama derecha descienden los otros dos subárboles, por lo que hay que hacer una llamada recursiva, obtener el lca entre las raíces de ellos y finalmente establecer las ramas hacia ellos.

El resultado de la unión se ve en la Figura 5.2. Los nodos que son el resultado de la obtención de lca , que es el glca , se muestran con un círculo grande, y las ramas que los unen con un trazo grueso. En este ejemplo, el glca coincide con los nodos que pertenecen a la unión y no a alguno de los subárboles. Esto no siempre es así, ya que es posible que un elemento del glca sea uno de los nodos de los cuales se obtiene el lca , tal como se verá en el ejemplo de la Figura 5.4. Esta unión tiene 4 hojas, una cantidad menor, 3, de nodos internos y la cantidad de nodos calculados es aún menor, 2. Con esta descripción alcanzó con la obtención de dos lca para obtener esos nodos. Con la versión en que se calcula el lca de cada par de hojas provenientes de subárboles diferentes hacen falta cinco cálculos.

En general, la cantidad de nodos del glca es obviamente no mayor que el de nodos internos de la unión, que a su vez, por tratarse de árboles compactos, es menor que la cantidad de hojas. Esto es

$$|\text{glca}(S)| < |\mathbf{L}(\text{union}(S))| = \left| \bigcup_{T \in S} \mathbf{L}(T) \right|. \quad (5.3)$$

Por lo tanto, y asumiendo que el lca se puede calcular en tiempo constante (lo cual es posible como se verá en el Apéndice A), tenemos que la cantidad de

hojas de los subárboles es una cota superior de la complejidad de la obtención del glca y la unión de subárboles.

5.2. Búsqueda de contextos representativos

Volvemos al problema planteado en esta tesis y nos proponemos encontrar los contextos representativos, que sabemos que son los bi-ramificables y los bi-sufijos, $R = \mathcal{B} \cup \mathcal{S}$. Además estableceremos una cota superior sobre la cantidad de esos contextos.

Los árboles de sufijos correspondientes a las cadenas de formación de contextos izquierda y derecha son, respectivamente, T_{y_ℓ} y T_{y_r} . Nuevamente identificamos a los nodos con los contextos que representan.

5.2.1. Obtención de los contextos representativos

El conjunto R es un subconjunto del producto cartesiano de $\mathbf{V}(T_{y_\ell})$ y $\mathbf{V}(T_{y_r})$. Considerando R como una relación, a cada elemento s_r de $\mathbf{V}(T_{y_r})$ le corresponde un subconjunto de $\mathbf{V}(T_{y_\ell})$ que denotamos $R_L(s_r)$. Es el conjunto de subcadenas de y_ℓ que forman contexto representativo con s_r .

Si encontramos los conjuntos $R_L(s_r)$ para todos los nodos s_r de T_{y_r} habremos encontrado todos los representativos. Sabemos que los bi-sufijos se obtienen trivialmente de una correspondencia uno a uno entre $\mathbf{L}(T_{y_\ell})$ y $\mathbf{L}(T_{y_r})$. Por otra parte el conjunto de los bi-ramificables, \mathcal{B} , es un subconjunto de $\mathbf{B}(T_{y_\ell}) \times \mathbf{B}(T_{y_r})$. Los vamos a obtener mediante el glca de subárboles compactos de T_{y_ℓ} . Los subárboles compactos los obtenemos recursivamente mediante union de otros subárboles compactos.

Recordemos que $T_{y_r}(s_r)$ es el subárbol de T_{y_r} con raíz en s_r . Llamamos $T_{y_\ell}^C(s_r)$ al conjunto de nodos de T_{y_ℓ} relacionados mediante $R_L(\cdot)$ con algún nodo de $T_{y_r}(s_r)$.

Teorema 5.17. *Para todo nodo s_r de T_{y_r} , $T_{y_\ell}^C(s_r)$ es un subárbol compacto de T_{y_ℓ} . Si, además, s_r es un nodo interno, llamamos S a la colección de subárboles compactos $T_{y_\ell}^C(\cdot)$ correspondientes a los hijos de s_r en T_{y_r} . Entonces, $R_L(s_r)$ es $\text{glca}(S)$ y $T_{y_\ell}^C(s_r)$ es igual a $\text{union}(S)$.*

Demostración. La demostración se hace por inducción en la altura de $T_{y_r}(s_r)$. Si la altura de $T_{y_r}(s_r)$ es 0, esto es, si s_r es una hoja, entonces tanto $R_L(s_r)$ como $T_{y_\ell}^C(s_r)$ son la hoja s_ℓ de T_{y_ℓ} , tal que (s_ℓ, s_r) es bi-sufijo, por lo que la proposición se cumple de manera trivial.

Si s_r es un nodo interno aceptamos como hipótesis inductiva que el teorema se cumple para todos los hijos de s_r , ya que las alturas de los árboles de los que

son raíces son menores que la altura de $T_{y_r}(\mathbf{s}_r)$. Supongamos que s_ℓ pertenece a $R_L(\mathbf{s}_r)$. Esto significa que (s_ℓ, \mathbf{s}_r) es bi-ramificable y esto, como consecuencia del Lema 4.20, a su vez implica que ocurren dos bi-sufijos $(s_\ell a \mathbf{u}_\ell, \mathbf{s}_r b \mathbf{u}_r)$, $(s_\ell c \mathbf{v}_\ell, \mathbf{s}_r d \mathbf{v}_r)$, con $a, b, c, d \in \mathcal{A}$, $a \neq c$, $b \neq d$. Por la hipótesis inductiva, el componente izquierdo del primero de estos bi-sufijos es una hoja en el subárbol compacto correspondiente al hijo de \mathbf{s}_r que desciende por la rama b (que es $R(\mathbf{s}_r, b)$, el representativo en el caso uni-direccional de \mathbf{s}_r, b). Esto es, $s_\ell a \mathbf{u}_\ell \in \mathbf{L}(T_{y_\ell}^C(R(\mathbf{s}_r, b)))$. Y por la misma razón $s_\ell c \mathbf{v}_\ell \in \mathbf{L}(T_{y_\ell}^C(R(\mathbf{s}_r, d)))$. Como $b \neq d$, $R(\mathbf{s}_r, b)$ y $R(\mathbf{s}_r, d)$ son hijos diferentes de \mathbf{s}_r , por lo que sus correspondientes árboles en S son diferentes. Entonces $s_\ell a \mathbf{u}_\ell$ y $s_\ell c \mathbf{v}_\ell$ son hojas de dos subárboles compactos diferentes de S , por lo que el lca de ambas, s_ℓ , pertenece al glca de los subárboles.

Recíprocamente, supongamos que s_ℓ pertenece al glca de los subárboles compactos correspondientes a hijos de \mathbf{s}_r , $T_{y_\ell}^C(R(\mathbf{s}_r, b))$ y $T_{y_\ell}^C(R(\mathbf{s}_r, d))$, con $b \neq d$. Esto significa que s_ℓ es el lca de un par de hojas, tomadas una de cada árbol. Sean estas hojas $s_\ell a \mathbf{u}_\ell$ y $s_\ell c \mathbf{v}_\ell$ respectivamente. Como s_ℓ es el lca de las hojas se cumple que $a \neq c$. Como $s_\ell a \mathbf{u}_\ell$ es una hoja de $T_{y_\ell}^C(R(\mathbf{s}_r, b))$ hay una hoja $\mathbf{s}_r b \mathbf{u}_r$ en $T_{y_r}(\mathbf{s}_r)$ tal que $(s_\ell a \mathbf{u}_\ell, \mathbf{s}_r b \mathbf{u}_r)$ es un bi-sufijo. De manera similar se llega a que $(s_\ell c \mathbf{v}_\ell, \mathbf{s}_r d \mathbf{v}_r)$ es un bi-sufijo. Por lo tanto, por Definición 4.19, (s_ℓ, \mathbf{s}_r) es bi-ramificable, por lo que $s_\ell \in R_L(\mathbf{s}_r)$.

Quedó demostrado que $R_L(\mathbf{s}_r)$ es $\text{glca}(S)$. Por definición, los nodos de $T_{y_\ell}^C(\mathbf{s}_r)$ son la unión de $R_L(\mathbf{s}_r)$ y los nodos de S . Por lo tanto, según la Ecuación (5.2), los nodos de $T_{y_\ell}^C(\mathbf{s}_r)$ son los mismos que los de $\text{union}(S)$. Al tener los mismos nodos, si a $T_{y_\ell}^C(\mathbf{s}_r)$ lo consideramos como un árbol compacto es el árbol soporte de las mismas hojas que $\text{union}(S)$ y por lo tanto es el mismo árbol. \square

El Algoritmo 3 permite obtener $T_{y_\ell}^C(\mathbf{s}_r)$.

Algoritmo 3 `subarbol_asociado`($\mathbf{s}_r : \mathbf{B}(T_{y_r})$)

Salida: $T_{y_\ell}^C(\mathbf{s}_r)$

$S \leftarrow \emptyset$

Para cada $\mathbf{t}_r \in \text{children}(T_{y_r}, \mathbf{s}_r)$

$$S \leftarrow S \cup \begin{cases} \mathbf{t}_\ell, \text{ tal que } (\mathbf{t}_\ell, \mathbf{t}_r) \text{ es bi-sufijo} & \text{si } \mathbf{t}_r \in \mathbf{L}(T_{y_r}) \\ \text{subarbol_asociado}(\mathbf{t}_r) & \text{si } \mathbf{t}_r \in \mathbf{B}(T_{y_r}) \end{cases}$$

devuelve $\text{union}(S)$

Como conocemos los subárboles compactos correspondientes a las hojas de T_{y_r} , si recorremos T_{y_r} en post-orden, y para cada nodo obtenemos la unión de los subárboles compactos correspondientes a sus hijos, al procesar la raíz de T_{y_r} habremos encontrado todos los contextos representativos. Y sabemos que para obtener la unión de subárboles alcanza con obtener el glca y conectar nodos mediante ramas.

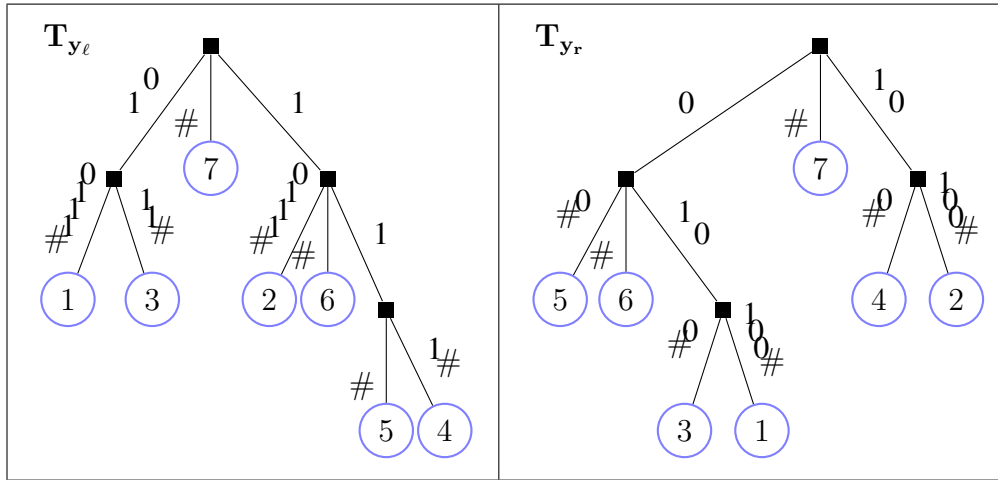


Figura 5.3: Árboles de las cadenas de formación de contextos de $\begin{pmatrix} y_\ell \\ y_r \end{pmatrix} = \begin{pmatrix} 010111\# \\ 010100\# \end{pmatrix}$.

Ejemplo 5.18. Seguimos con el ejemplo de la Sección 4.2 : $\mathbf{x} = \begin{pmatrix} 1110101 \\ 0010101 \end{pmatrix}$, de donde $\begin{pmatrix} y_\ell \\ y_r \end{pmatrix} = \begin{pmatrix} 010111\# \\ 010100\# \end{pmatrix}$, cuyos árboles de sufijos correspondientes repetimos en la Figura 5.3.

Para cada hoja i , $1 \leq i \leq 7$ de T_{y_r} , el subárbol compacto de T_{y_ℓ} correspondiente es la hoja i de T_{y_ℓ} . El subárbol que corresponde al nodo 010 de T_{y_r} se obtiene mediante la unión de los subárboles correspondientes a las hojas 1 y 3, que son las hojas 1 y 3 respectivamente de T_{y_ℓ} . Esta unión tiene 01 como raíz y sus dos hijos son las hojas mencionadas. El glca es sólo la raíz, 01. Entonces $(01, 010)$ es el único contexto bi-ramificable cuyo contexto derecho es 010.

Después se procesa el nodo 0 de T_{y_r} . Esto es lo que se vio en la Figura 5.2. Hay que notar que en esa figura la estructura del árbol base es la de T_{y_ℓ} , y que los tres subárboles de arriba son los correspondientes a los hijos de 0 en T_{y_r} . El glca está formada por dos nodos, entonces, por el Teorema 5.17, $R_L(0)$ es $\{\epsilon, 1\}$. Por lo tanto $(\epsilon, 0)$ y $(1, 0)$ también son bi-ramificables.

Con el nodo 10 de T_{y_r} se obtiene un subárbol de T_{y_ℓ} con raíz en 1 cuyos hijos son las hojas 2 y 4. El contexto $(1, 10)$ es el único bi-ramificable.

Finalmente se procesa la raíz de T_{y_r} . Esto se ve en la Figura 5.4.

En ese ejemplo se ve que dos de los tres nodos que se calcularon pertenecen a alguno o algunos de los subárboles que constituyen el parámetro.

En la parte de arriba están los subárboles correspondientes a los tres hijos de ϵ

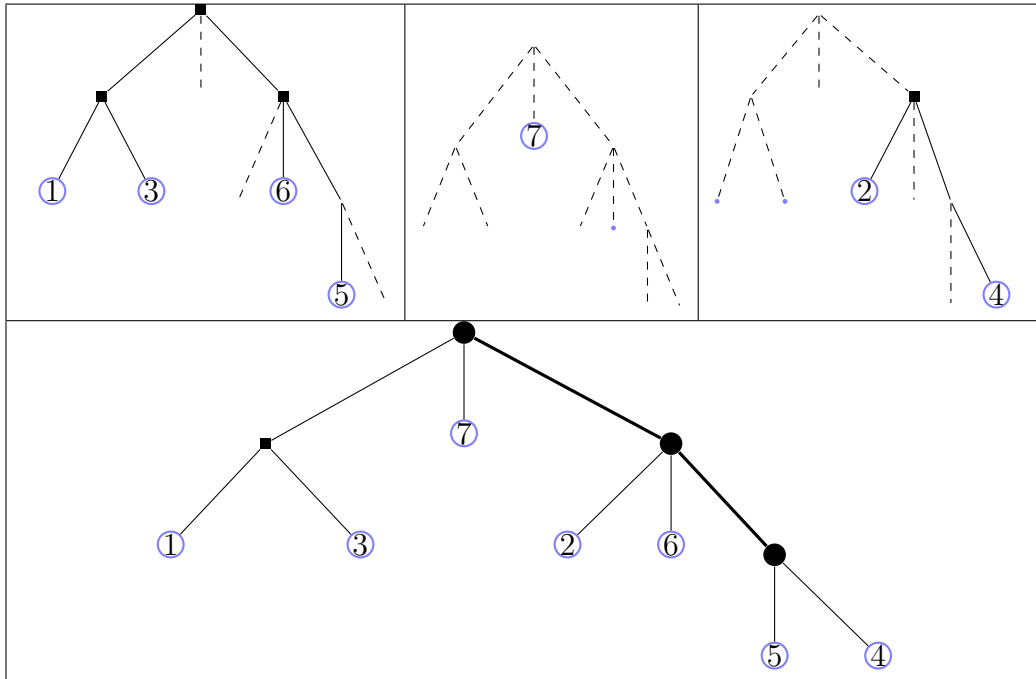


Figura 5.4: Subárbol asociado a la raíz

en T_{y_r} . El lca de las raíces de los subárboles es ϵ , por lo que este nodo pertenece al glca , lo cual se indica en la figura con el círculo. Por la rama 0, sólo el primer subárbol tiene descendientes, por lo tanto alcanza con establecer una rama hacia 01, que es el hijo por esa rama. El subárbol del medio es descendiente por la rama # y los otros dos no tienen descendencia por esa rama. Finalmente el primer y tercer subárbol son los que tienen nodos que descienden de la raíz obtenida por la rama 1. Para el primer subárbol los descendientes forman el árbol soporte de las hojas 6 y 5, mientras que al tercer subárbol corresponde el soporte de las hojas 2 y 4. El lca de las raíces de estos dos subárboles es el nodo 1, que por lo tanto, pertenece al glca . De este nodo descienden por la rama 0 sólo la hoja 2 del tercer subárbol y por la rama # sólo la hoja 6 del primer subárbol. Pero por la rama 1 descienden hojas en ambos subárboles, la 5 del primero y la 4 del tercero. Por lo tanto se debe obtener el lca de ambas hojas que es 11. Entonces se completó el conjunto \mathcal{B} de bi-ramificables con (ϵ, ϵ) , $(1, \epsilon)$ y $(11, \epsilon)$.

5.2.2. Cantidad de contextos representativos

Ahora vamos a establecer una cota acerca de la cantidad de contextos representativos.

Denotamos por $h(\cdot)$ la cantidad de niveles de un árbol.

Veremos que la cantidad de contextos representativos está acotada superiormente por $n \times h(T_{y_r})$ porque en cada nivel de T_{y_r} hay a lo sumo n contextos representativos.

Lema 5.19. *La cantidad de contextos representativos encontrados con el método de la Sección 5.2.1 está acotada superiormente por $n \times h(T_{y_r})$.*

Demostración. Notemos que por cada nodo s_r de T_{y_r} se cumple que la cantidad de contextos representativos está acotado por la cantidad de hojas que en T_{y_r} descienden de s_r , esto es

$$|R_L(s_r)| < |\mathbf{L}(T_{y_r}(s_r))|,$$

lo que se deduce de la Ecuación (5.3) ya que en el Teorema 5.17 hemos probado que $R_L(s_r)$ es el glca del conjunto S de subárboles $T_{y_\ell}^C(\cdot)$ correspondientes a los hijos de s_r , y sabemos que las hojas que descienden de s_r están en correspondencia uno a uno con las de los subárboles en S .

Como los conjuntos de hojas que descienden de nodos del mismo nivel de T_{y_r} son disjuntos, y el número de hojas en T_{y_r} es el largo n de la cadena, llegamos a que n es una cota superior ajustada de la cantidad de hojas que descienden de todos los nodos del nivel.

Se concluye lo propuesto al tratar cada nivel de T_{y_r} . □

Finalmente ajustamos la cota de cantidad de contextos representativos.

Teorema 5.20. *La cantidad de contextos representativos correspondientes a la secuencia \mathbf{x} está acotada superiormente por $n \cdot \min \{h(T_{y_\ell}), h(T_{y_r})\}$, donde T_{y_ℓ} y T_{y_r} son los árboles de sufijos de las cadenas de formación de contextos izquierda y derecha respectivamente correspondientes a \mathbf{x} .*

Demostración. En esta sección hemos descrito la obtención de los contextos representativos mediante una recorrida de T_{y_r} . Podríamos haber hecho algo similar recorriendo T_{y_ℓ} . Para esta recorrida se puede establecer un postulado análogo al Lema 5.19. Se obtienen entonces dos expresiones que son cota superior de la cantidad de contextos representativas, lo cual prueba el enunciado del teorema. □

En la Sección 6.4 veremos que construir un grafo con todos los contextos representativos es lo que determina la complejidad de la obtención de un conjunto óptimo de contextos, y que cada contexto representativo se obtiene en tiempo constante. Por lo tanto el resultado que aquí se obtuvo determina la complejidad del algoritmo completo.

>

Capítulo 6

Grafos Compactos de Contextos Bi–direccionales

En la Sección 3.2 se estableció un algoritmo que permite obtener un conjunto óptimo de contextos bi–direccionales para el caso de longitudes acotadas. Se hacía notar que se obtenía un árbol que de manera implícita recorría un grafo, denominado grafo de contextos bi–direccionales. En este capítulo se describe el grafo, junto con su construcción y podado, correspondiente al caso de longitudes no acotadas.

Aquel grafo tenía un vértice para cada contexto posible, incluso los que no ocurrían. Del Corolario 4.30 sabemos que podemos restringir la atención a contextos bi–ramificables y a bi–sufijos. En este capítulo definimos una clase de grafos, los grafos compactos de contextos bi–direccionales cuyos vértices corresponden a esos contextos. Entonces, el problema en cadenas de símbolos de encontrar los contextos representativos y un conjunto óptimo de contextos es trasladado al problema de construir un grafo y podarlo. El resultado del podado es un árbol compacto cuyas hojas forman un conjunto óptimo de contextos bi–direccionales.

La construcción del grafo se basa en la obtención de los contextos representativos vista en la Sección 5.2.1. En esa sección se vio que aplicando el Algoritmo 3 a cada nodo s de T_{y_r} se obtiene un conjunto de nodos de T_{y_ℓ} que forman un subárbol compacto asociado a s . Lo que falta para completar el grafo es establecer las aristas entre los nodos de los subárboles asociados a cada nodo de T_{y_r} . Para ello, se verá el Algoritmo 4 que es una generalización del Algoritmo 3 mediante el cual para cada nodo de T_{y_r} se obtiene un grafo en lugar de un árbol.

En la Sección 6.2 se describe formalmente la construcción de un grafo. El podado se realiza junto con esa construcción como se explica en la Sección 6.3, en donde también se demuestra que se obtiene un conjunto óptimo de contextos bi–direccionales. Se estudia la complejidad en la Sección 6.4. En la Sección 6.5 se describen resultados experimentales.

6.1. Definición

Un *grafo compacto de contextos bi-direccionales* (o CBDCG por su nombre en inglés, “compact bi-directional context graph”) es un grafo enraizado, dirigido y acíclico, definido para cada secuencia x , de tal manera que hay una correspondencia uno a uno entre los vértices del grafo y los contextos representativos de x . Con un abuso de notación a los vértices sumideros los llamaremos ‘hojas’, a los que tienen aristas salientes ‘internos’ y a los destinos de las aristas salientes, ‘hijos’. La raíz está etiquetada con (ϵ, ϵ) , los vértices internos corresponden a los bi-ramificables y las hojas representan los bi-sufijos. Cada vértice interno tiene dos conjuntos de hijos, ambos con más de un elemento. Para uno de los conjuntos, el de los *hijos izquierdos* del vértice, se establece una correspondencia uno a uno con las extensiones mínimas por izquierda que ocurren del vértice. Más precisamente, el hijo es el contexto representativo de la extensión mínima que le corresponde. O sea, si (s_ℓ, s_r) es un vértice interno y $(s_\ell a, s_r)$ es una extensión que ocurre, entonces el hijo de (s_ℓ, s_r) que corresponde a $(s_\ell a, s_r)$ es $R(s_\ell a, s_r)$. Obviamente, cada uno de los componentes izquierdos de las aristas que salen de (s_ℓ, s_r) hacia sus hijos izquierdos comienzan con un símbolo diferente. Está claro que el conjunto de hijos izquierdos de (s_ℓ, s_r) es un conjunto válido relativo a (s_ℓ, s_r) , precisamente el conjunto válido mínimo por izquierda, definido en la Expresión (4.3). Consideraciones análogas valen para el conjunto de *hijos derechos* del vértice. Las aristas están etiquetadas con un par ordenado de cadenas sobre \mathcal{A}^* y al menos uno de sus componentes no es ϵ . Entonces, un CBDCG puede ser considerado como una generalización de los árboles compactos. Notemos, sin embargo que un vértice de un CBDCG puede ser hijo derecho de más de un vértice, y al mismo tiempo ser hijo izquierdo de los mismos u otros vértices.

En la Figura 6.1 se ve el grafo correspondiente al Ejemplo 5.18. Por ejemplo $(1, 10)$ es hijo derecho tanto de (ϵ, ϵ) como de $(1, \epsilon)$. El bi-sufijo 6 es hijo derecho de $(\epsilon, 0)$ e hijo derecho e izquierdo de $(1, 0)$.

Ejemplo 6.1. Supongamos que (s_ℓ, s_r) , $(s_\ell a, s_r)$ y $(s_\ell a, s_r c)$ son contextos bi-ramificables, y por lo tanto vértices del grafo. El vértice $(s_\ell a, s_r c)$ es hijo derecho de $(s_\ell a, s_r)$ y también sería hijo derecho de (s_ℓ, s_r) si no ocurriera $(s_\ell b, s_r c)$ con $b \neq a$. Vamos a ver un ejemplo para la aplicación eliminación de ruido en el que sucede esto (esta aplicación limita la libertad en la elección de las cadenas formadoras de contextos y_ℓ e y_r , y es por lo tanto un ejemplo más fuerte). Para que el ejemplo sea más sencillo usamos el alfabeto $\{0, 1\}$, tomamos $(s_\ell, s_r) = (\epsilon, \epsilon)$ y no tendremos en cuenta la primera y la última posición para no necesitar el símbolo marcador. Sea $x = 00000101011$. Comprobamos que

- (ϵ, ϵ) es bi-ramificable porque ocurren $(0, 0)$ (posiciones 2, 3, 4 y 8) y $(1, 1)$ (posiciones 7 y 9).

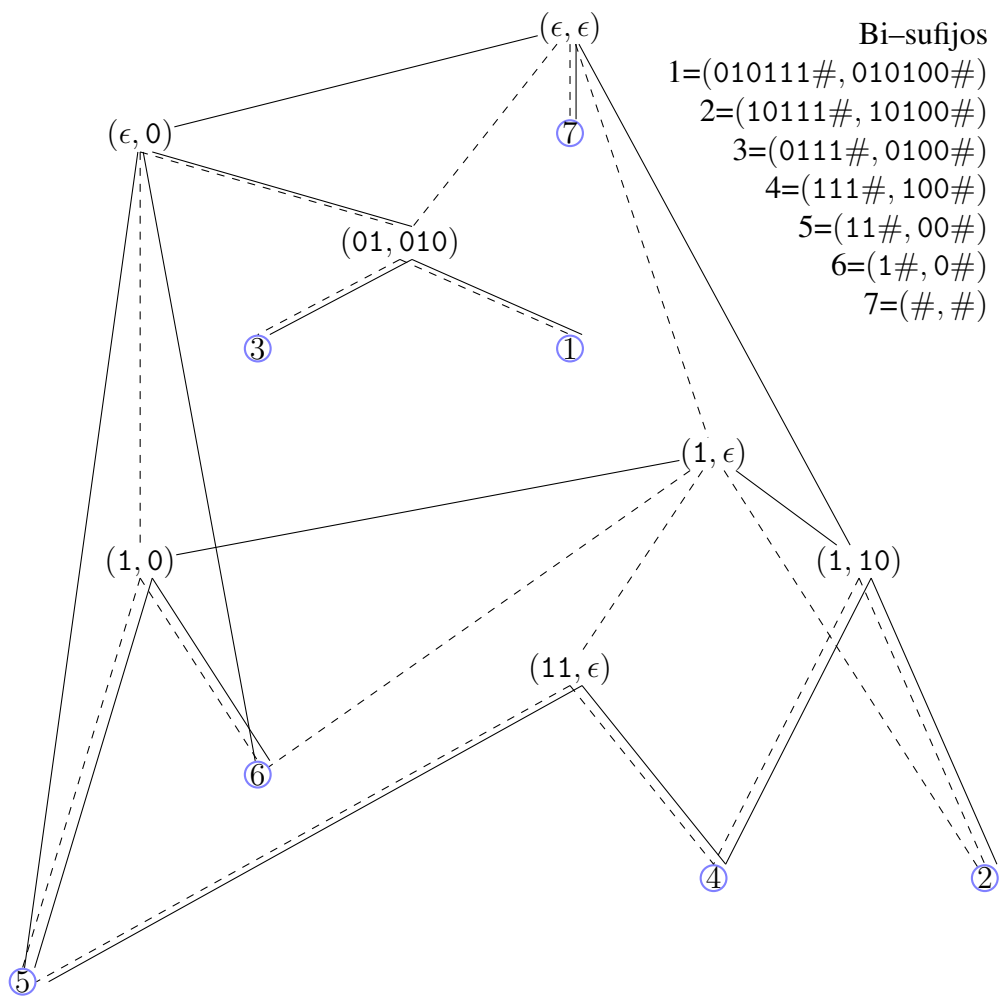


Figura 6.1: Grafo compacto de contextos bi-direccionales para Ejemplo 5.18. Se muestran todos los vertices del grafo. Los bi-sufijos se representan con el numero de hoja y los bi-ramificables con su contexto. Las aristas de trazo continuo muestran las relaciones padre-hijo derecho, ordenadas para cada vertice de izquierda a derecha segun las extensiones por la arista 0, # y 1 en caso de que las haya. Las aristas de trazo discontinuo muestran las relaciones padre-hijo izquierdo. Todas las aristas estan orientadas desde arriba hacia abajo.

- $(0, \epsilon)$ es bi-ramificable porque ocurren $(00, 0)$ (posiciones 3 y 4) y $(01, 1)$ (posición 10).
- $(0, 0)$ es bi-ramificable porque ocurren $(00, 00)$ (posición 3) y $(01, 01)$ (posición 8).
- $(1, 0)$ no ocurre.

Por lo tanto $(0, 0)$ es hijo derecho de (ϵ, ϵ) y de $(0, \epsilon)$.

Un CBDCG se construye para obtener un conjunto óptimo definido por las hojas de un árbol óptimo, el cual se obtiene podando el grafo. Por eso, de cada vértice del grafo, además de la información de su estructura, debe mantenerse otra que permita determinar de qué manera se obtiene un conjunto óptimo relativo a él. Para cada contexto (s_ℓ, s_r) esto se conoce con $\text{split}(s_\ell, s_r)$, que indica si el conjunto óptimo se obtiene mediante la partición izquierda, la derecha, o con la ‘no partición’, o sea, si el vértice pertenece al conjunto óptimo. Para poder establecer este valor hay que conocer el del peso óptimo $w_{opt}(s_\ell, s_r)$, que a su vez depende del vector de conteos $c(s_\ell, s_r)$ vía el peso del vértice, $w(s_\ell, s_r)$. En la Sección 3.2 se presentó el Algoritmo 1 que encuentra el conjunto óptimo para el caso de contextos de longitud acotada. El correspondiente al caso actual se presentará en el Algoritmo 6 que usa información que se va estableciendo durante la construcción del grafo.

Con una analogía a los subárboles, definimos el *subgrafo enraizado* correspondiente a un vértice de un CBDCG como el conjunto de vértices que descienden del vértice considerado y las aristas que los unen. En adelante, siempre que se hable de subgrafo significa subgrafo enraizado. El subgrafo correspondiente a (s_ℓ, s_r) se denota $G(s_\ell, s_r)$. El conjunto de todos los subgrafos es \mathcal{G} . Está implícito que son los subgrafos generados por la secuencia \mathbf{x} .

Entonces damos la especificación para cada subgrafo $G \in \mathcal{G}$.

- $\text{root}(G) = (\text{left}(G), \text{right}(G)) \in R$ es el contexto representativo que corresponde a la raíz del subgrafo.

En lo que queda de esta especificación $(s_\ell, s_r) = (\text{left}(G), \text{right}(G))$.

- Si $\text{root}(G) \in \mathcal{B}$, entonces $G^L, G^R \subset \mathcal{G}$ son los conjuntos de subgrafos cuyas raíces son los hijos izquierdos y derechos respectivamente de la raíz de

G . Cada conjunto tiene al menos dos elementos. Si $(s_\ell a, s_r)$ (resp. $(s_\ell, s_r a)$) ocurre, entonces $G^L[a] = G(R(s_\ell a, s_r))$ (resp. $G^R[a] = G(R(s_\ell, s_r a))$).

Si $\text{root}(G) \in \mathcal{S}$ (conjunto de bi-sufijos), entonces $G^L = \emptyset$ y $G^R = \emptyset$.

- $c(s_\ell, s_r)$, $w_{opt}(s_\ell, s_r)$, $\text{split}(s_\ell, s_r)$ son el vector de conteos, el peso óptimo y el split respectivamente de la raíz de G . Se definen como en las Ecuaciones (4.5), (4.6) y (4.7) respectivamente.

6.2. Construcción del CBDCG

En la Sección 5.2 se describió cómo calcular los contextos representativos y también se obtuvo una cota superior de la cantidad de dichos contextos. Si sólo interesara obtener esos contextos el Algoritmo 3 sería suficiente. Pero como el objetivo es construir el CBDCG, además de obtener los contextos hay que establecer la estructura del grafo, esto es, las relaciones padre-hijo, y asignar los valores del vector de conteos, peso óptimo e indicador de split.

Sin embargo, del Algoritmo 3 también se puede deducir parte de la estructura del grafo. Para cada nodo s_r de T_{y_r} se obtuvo $R_L(s_r)$, el conjunto de nodos de T_{y_ℓ} con los que forma contexto representativo. Además se conoce el conjunto $T_{y_\ell}^C(s_r)$ que contiene a $R_L(s_r)$ y a todos los nodos de T_{y_ℓ} con los que los descendientes de s_r en T_{y_r} forman contexto representativo. Este conjunto $T_{y_\ell}^C(s_r)$ tiene estructura de árbol (Teorema 5.17), de la cual se pueden obtener las aristas padre-hijo izquierdo en el grafo que se quiere construir. En esta sección, extendiendo este algoritmo con el Algoritmo 4, se va a completar la construcción del grafo.

La construcción del CBDCG empieza por obtener las cadenas de formación de contextos y sus árboles de sufijos, T_{y_ℓ} y T_{y_r} . Entonces se procede a un recorrido en post-orden de uno de ellos. Por razones de complejidad demostradas en el Teorema 5.20, y en las que se insistirá en la Sección 6.4, se elige aquel cuya altura sea menor. Sin pérdida de generalidad, la siguiente descripción supone que ese es el derecho, T_{y_r} . El objetivo es encontrar para cada s_r en $V(T_{y_r})$ el conjunto de nodos de T_{y_ℓ} que forman un contexto representativo con s_r , o sea $R_L(s_r)$, establecer la estructura, esto es, las relaciones padre-hijo, y asignar los valores del vector de conteos, peso óptimo e indicador de split.

Para cada nodo s_r de T_{y_r} , obtendremos un subgrafo, $g(s_r)$, cuyos vértices están en correspondencia con el conjunto $T_{y_\ell}^C(s_r)$, o sea con los componentes izquierdos de los contextos representativos cuyos componentes derechos son las extensiones de s_r . Esto significa que es el subgrafo $G(s_\ell, s_r)$, siendo $(s_\ell, s_r) = R(\epsilon, s_r)$.

En el paso inicial del recorrido en post-orden se relaciona cada hoja de T_{y_r} con la hoja correspondiente de T_{y_ℓ} , de tal manera que se crea un vértice del CBDCG que corresponde a un bi-sufijo. El subgrafo que corresponde a la hoja de T_{y_r} tiene

un único vértice. El resto de la información del vértice, esto es vector de conteos, peso óptimo y split se obtiene de manera trivial.

En el paso general, al procesar el nodo interno s_r de T_{y_r} , el algoritmo usa T_{y_ℓ} para obtener el conjunto $R_L(s_r)$ formado por los contextos izquierdos que forman un contexto representativo con s_r . Por cada elemento de $R_L(s_r)$ se obtiene un vértice del CBDCG con s_r como componente derecho. Además se establece la estructura del grafo.

El recorrido de T_{y_r} se hace con la función recursiva **graph**, descrita en el Algoritmo 4.

En ese algoritmo, s_r es un nodo interno de T_{y_r} y recordamos que **children** devuelve el conjunto de hijos del nodo pasado como parámetro. La variable S es un conjunto de subgrafos.

Se puede ver que este algoritmo es una analogía del Algoritmo 3 en la que se devuelven grafos en lugar de árboles.

Algoritmo 4 $\text{graph}(s_r : \mathbf{B}(T_{y_r}))$

Salida: $g(s_r) : \mathcal{G}$

Var $S : \text{set}(\mathcal{G})$

$S \leftarrow \emptyset$

Para cada $t_r \in \text{children}(T_{y_r}, s_r)$

$$S \leftarrow S \cup \begin{cases} g(t_r) & \text{si } t_r \in \mathbf{L}(T_{y_r}) \\ \text{graph}(t_r) & \text{si } t_r \in \mathbf{B}(T_{y_r}) \end{cases}$$

$g(s_r) \leftarrow \text{union}(S, \epsilon, s_r)$

devuelve $g(s_r)$

Postcondición: $\text{root}(g(s_r)) = R(\epsilon, s_r)$

La primera invocación de **graph** se hace desde el algoritmo principal (Algoritmo 7) pasándole ϵ como parámetro. Después del recorrido en post-orden, finaliza al procesar la raíz de T_{y_r} y la salida es el CBDCG.

Como se ve en el Algoritmo 4, $g(s_r)$ es la salida de la función **union**. Esta función es la más importante de la construcción y se describe en el Algoritmo 5. De acuerdo a la especificación, **root** es la raíz de un subgrafo, con **left** y **right** se representan los componentes izquierdo y derecho respectivamente de la raíz del subgrafo, el conjunto de grafos cuyas raíces son los hijos izquierdos (resp. derechos) de la raíz de un subgrafo G se denota G^L (resp. G^R), y $G^L[a] \in G^L$ tiene como raíz el hijo izquierdo de **root**(G) que extiende **left**(G) según la rama a . O sea $\text{root}(G^L[a]) = R(\text{left}(G)a, \text{right}(G))$. La función **lca** obtiene el ancestro común más profundo de un conjunto de nodos del árbol T_{y_ℓ} , que es el prefijo común más largo de los contextos (izquierdos) representados por esos nodos. Como se había adelantado, este algoritmo es una generalización del Algoritmo 2.

Algoritmo 5 $union(S : \text{set}(\mathcal{G}), \mathbf{v}_\ell : \mathbf{B}(T_{y_\ell}), \mathbf{s}_r : \mathbf{B}(T_{y_r}))$

- 1: **Salida:** $U : \mathcal{G}$
 - 2: **Var** $S' : \text{set}(\mathcal{G})$
 - 3: **Precondición 1:** $|S| \geq 2$
 - 4: **Precondición 2:** Las raíces de los grafos de S forman el conjunto válido mínimo por derecha (Expresión 4.4) de $(\mathbf{v}_\ell, \mathbf{s}_r)$:
 $\{\text{root}(G) : G \in S\} = \{\mathbf{R}(\mathbf{v}_\ell, \mathbf{s}_r b) : b \in \mathcal{A}, \mathcal{I}(\mathbf{v}_\ell, \mathbf{s}_r b) \neq \emptyset\}$
 - 5: $\mathbf{s}_\ell = \text{lca}\{\text{left}(G) : G \in S\}, \mathbf{s}_\ell \in \mathbf{B}(T_{y_\ell})$
 - 6: $\text{root}(U) \leftarrow (\mathbf{s}_\ell, \mathbf{s}_r)$
 - 7: $U^R \leftarrow S$
 - 8: **Para cada** $\mathbf{t}_\ell \in \text{children}(T_{y_\ell}, \mathbf{s}_\ell)$
 - 9: Sea $a \in \mathcal{A}$ tal que $a = [\mathbf{t}_\ell]_{|\mathbf{s}_\ell|+1}$, o sea $\mathbf{s}_\ell a \preceq \mathbf{t}_\ell$
 - 10: $S' \leftarrow \emptyset$
 - 11: **Para cada** $G \in S$
 - 12: $S' \leftarrow S' \cup \begin{cases} G & \text{si } \mathbf{t}_\ell \preceq \text{left}(G) \\ G^L[a] & \text{si } \mathbf{s}_\ell = \text{left}(G) \text{ y } \text{root}(G^L[a]) \text{ ocurre} \\ \emptyset & \text{en otro caso.} \end{cases}$
 - 13: **Si** $S' \neq \emptyset$ **entonces**
 - 14: $U^L[a] \leftarrow \begin{cases} G' & \text{si } S' = \{G'\} \\ \text{union}(S', \mathbf{t}_\ell, \mathbf{s}_r) & \text{si } |S'| \geq 2. \end{cases}$
 - 15: {actualizar vector de conteos, peso óptimo, split de $\text{root}(U)$ }
 - 16: **devuelve** U
 - 17: **Postcondición:** $\text{root}(U) = \mathbf{R}(\mathbf{v}_\ell, \mathbf{s}_r)$
-

Para cada nodo $s_r \in \mathbf{B}(T_{y_r})$ la llamada inicial de la función **union** invocada desde el Algoritmo 4 pasa como primer parámetro el conjunto de los subgrafos $g(t_r)$ para todos los hijos t_r de s_r . La salida de la llamada inicial de **union** es el subgrafo $g(s_r)$ cuya raíz es $(s_\ell, s_r) = R(\epsilon, s_r)$. Mientras se genera $g(s_r)$, **union** obtiene todos los vértices correspondientes a los contextos representativos cuyo componente derecho es s_r . Cada llamada recursiva a **union** obtiene un vértice que es la raíz del subgrafo U que devuelve. Para eso en la línea 5 calcula s_ℓ , que es el lca de los componentes izquierdos de las raíces de los grafos del conjunto parámetro S . Estas raíces son los hijos derechos del vértice que se va a obtener, por lo que inmediatamente se puede establecer en la línea 7 las aristas padre–hijo derecho. Estas aristas están determinadas por la correspondencia establecida en la Precondición 2, es decir, si $\text{root}(G) = R(v_\ell, s_r b)$ entonces $U^R[b] = G$. Para establecer las relaciones padre–hijo izquierdo se aplica el bucle de la línea 8, que es el análogo del de la línea 4 del Algoritmo 2. Tal como en ese algoritmo, es posible que deban hacerse llamadas recursivas a **union**. Para cada hijo t_ℓ de s_ℓ en T_{y_ℓ} se considera el símbolo a que es el primer símbolo de u , donde $t_\ell = s_\ell u$. Es decir el símbolo a identifica la rama por la cual desciende t_ℓ . En el bucle de la línea 11 se obtiene el subgrafo de cada G en S que puede contribuir a la formación de $U^L[a]$, el subgrafo izquierdo de $\text{root}(U)$ según la rama a . Es el subgrafo cuyos vértices son extensiones de $(t_\ell, \text{right}(G))$. Puede no existir, ser el propio G , o ser $G^L[a]$. Se lo incluye en un conjunto S' . Si S' no es vacío entonces en la línea 14 se establece la relación padre–hijo izquierdo entre $\text{root}(U)$ y $\text{root}(U^L[a])$. El subgrafo $U^L[a]$ puede ser un subgrafo que ya existía si S' tiene un único elemento, o el resultado de una llamada recursiva si tiene más de uno. Finalmente, en la línea 15 se actualiza la información del vértice, esto es, vectores de conteos, peso óptimo y **split**, que serán utilizada para el podado del grafo como se verá en la Sección 6.3. Para consideraciones de complejidad computacional notemos que el número total de llamadas a **union** es la cantidad de vértices internos del grafo.

Veamos que las precondiciones del Algoritmo 5 se cumplen. El algoritmo es invocado desde dos lugares, la llamada inicial desde el Algoritmo 4 y la recursiva desde la línea 14.

En la llamada hecha desde el Algoritmo 4 la Precondición 1 se cumple porque s_r es un nodo interno de T_{y_r} y por lo tanto S tiene más de un elemento. Para cada símbolo b tal que $(\epsilon, s_r b)$ ocurre hay en T_{y_r} exactamente un t_r hijo de s_r , y por lo tanto un grafo en S . La raíz de $g(t_r)$ es $R(\epsilon, t_r)$ si t_r es una hoja porque es el bi–sufijo cuyo componente derecho es t_r . Y si t_r no es hoja, por la postcondición del algoritmo, la raíz de $g(t_r)$ también es $R(\epsilon, t_r)$. Notando que el valor de v_ℓ en la llamada a **union** es ϵ , se concluye que se cumple la Precondición 2.

Para considerar las llamadas recursivas vamos a establecer dos lemas.

Lema 6.2. *Si (t_ℓ, t_r) es extensión estricta de (s_ℓ, s_r) y prefijo del representativo*

de (s_ℓ, s_r) entonces ambos tienen el mismo representativo.

$$\text{Si } (s_\ell, s_r) \prec (t_\ell, t_r) \preceq R(s_\ell, s_r) \text{ entonces } R(t_\ell, t_r) = R(s_\ell, s_r).$$

Demostración. Por las relaciones de prefijo dadas en las hipótesis se sabe que (t_ℓ, t_r) pertenece a la misma clase que (s_ℓ, s_r) . Por lo tanto sus representativos son el mismo. \square

Lema 6.3. Si $R(s_\ell, s_r) = (s_\ell, t_r)$ entonces $R(s_\ell a, s_r) = R(s_\ell a, t_r)$.

Demostración. Si $s_r = t_r$ la proposición se cumple de manera trivial.

Supongamos, en otro caso, que $s_r b \preceq t_r$. La hipótesis implica que no ocurre $(s_\ell, s_r d)$, con $d \neq b$, porque si ocurriera sería $R(s_\ell, s_r) = (s_\ell, s_r)$. Tenemos $(s_\ell a, s_r) \prec (s_\ell a, t_r) \preceq R(s_\ell a, t_r)$. Sea R' el representativo de $(s_\ell a, s_r)$. No se puede cumplir $R(s_\ell a, t_r) \prec R'$ porque el representativo de $(s_\ell a, s_r)$ sería $R(s_\ell a, t_r)$ y no R' . Tampoco se puede cumplir $(s_\ell a, t_r) \preceq R' \prec R(s_\ell a, t_r)$ porque el representativo de $(s_\ell a, t_r)$ sería R' y no $R(s_\ell a, t_r)$. Supongamos, para llegar a una contradicción, que se cumple $(s_\ell a, s_r) \preceq R' \prec (s_\ell a, t_r)$ y por lo tanto $R' = (s_\ell a, s_r v_r)$ para alguna cadena v_r , con $s_r v_r \prec t_r$. Recordando que $s_r b \preceq t_r$ y como R' es bi-ramificable tiene que existir $d \neq b$ tal que ocurra $(s_\ell a, s_r d)$ y por lo tanto también ocurre $(s_\ell, s_r d)$ en contradicción con lo que se había establecido al principio. La última suposición acerca de R' fue lo que llevó a esta contradicción. La única opción que queda es $R' = R(s_\ell a, t_r)$ \square

En las llamadas recursivas la Precondición 1 se cumple de manera explícita en la línea 14. Notemos que por la Precondición 2, se cumple que $v_\ell \preceq \text{left}(G)$ para todo $G \in S$, de donde $v_\ell \preceq s_\ell$, por lo que para todo t_ℓ del bucle de la línea 8 se cumple $v_\ell \prec t_\ell$. Por lo tanto en el bucle de la línea 11 si $(t_\ell, s_r b)$ ocurre existe un único grafo G para el que se cumple

$$\text{root}(G) = (\text{left}(G), \text{right}(G)) = R(v_\ell, s_r b) \preceq R(t_\ell, s_r b). \quad (6.1)$$

Si en la línea 12 se da el primer caso, $t_\ell \preceq \text{left}(G)$, y como ya se sabía que $s_r b \preceq \text{right}(G)$, tenemos $(t_\ell, s_r b) \preceq (\text{left}(G), \text{right}(G))$ de donde, junto con (6.1) y el Lema 6.2, tenemos $\text{root}(G) = R(t_\ell, s_r b)$. Y en este caso, G es el grafo que se incluye en S' .

Para el segundo caso de la línea 12, notemos que de $v_\ell \preceq s_\ell$, por el Lema 6.2 y por la Precondición 2 se cumple $R(s_\ell, s_r b) = \text{root}(G)$. Como en este caso $s_\ell = \text{left}(G)$ estamos en las hipótesis del Lema 6.3. Por lo tanto $R(s_\ell a, s_r b) = R(\text{left}(G)a, \text{right}(G))$, que es $\text{root}(G^L[a])$. Y $G^L[a]$ es el grafo que se incluye en S' .

Por lo tanto en S' se incluyen los grafos cuyas raíces son los representativos de $(t_\ell, s_r b)$ para todos los $(t_\ell, s_r b)$ que ocurren. Esas raíces son el conjunto válido mínimo por derecha de (t_ℓ, s_r) cuyos componentes son los parámetros de

union de la línea 14. Queda demostrado que se cumplen las precondiciones del Algoritmo 5.

Entonces, el resultado de *union* es el subgrafo constituido por todos los contextos representativos que son extensión de $(\mathbf{v}_\ell, \mathbf{s}_r)$. Como en el Algoritmo 4 esta función se invoca con $\mathbf{v}_\ell = \epsilon$ se concluye el siguiente teorema:

Teorema 6.4. *El algoritmo $\mathbf{graph}(\mathbf{s}_r)$ construye correctamente la estructura del subgrafo que contiene todos los vértices cuyo componente derecho es una extensión de \mathbf{s}_r .*

Decir que lo construye correctamente significa que hay una correspondencia uno a uno entre los vértices del subgrafo y los contextos representativos cuya componente derecha es extensión de \mathbf{s}_r , y que los hijos de los vértices internos son los representativos de las extensiones mínimas que ocurren.

6.3. Podado

El podado del grafo se hace al mismo tiempo que se lo construye cuando en la función *union* se asigna el valor del indicador de *split* al vértice creado. Por lo tanto está implícitamente definido el árbol óptimo que tiene raíz en ese vértice.

En un paso final se obtiene formalmente el árbol óptimo que corresponde a \mathbf{x} . Se hace recursivamente, comenzando en la raíz usando la función *prune* descrita en el Algoritmo 6.

Algoritmo 6 *prune* ($G : \mathcal{G}$)

Salida: $T_{opt} : \mathcal{T}_{opt}(\mathbf{root}(G))$

Sea $\mathbf{split} = \mathbf{split}(\mathbf{root}(G))$

$$T_{opt} = \{\mathbf{root}(G)\} \cup \begin{cases} \emptyset & \text{si } \mathbf{split} = 'N' \\ \bigcup_{G' \in G^L} \mathbf{prune}(G') & \text{si } \mathbf{split} = 'L' \\ \bigcup_{G' \in G^R} \mathbf{prune}(G') & \text{si } \mathbf{split} = 'R' \end{cases}$$

devuelve T_{opt}

Para ver que el valor de *split* se asigna correctamente describimos ahora la línea 15 del final del Algoritmo 5. Por el Teorema 6.4 se sabe que para cada vértice del CBDCG cada hijo es asignado correctamente de acuerdo a la definición de CBDCG. Y los conjuntos de hijos así definidos son los conjuntos en que se basan las Ecuaciones (4.5), (4.6) y (4.7). En el Algoritmo 5 se asignan los valores

a la raíz de U usando estas ecuaciones que sabemos que son correctas. Por lo tanto $\text{split}(U)$ establece la partición correcta.

Queda por ver que la estructura devuelta es un árbol. Lo es de manera trivial si $\text{split} = 'N'$ porque tiene un único nodo. En otro caso a un conjunto de árboles se le agrega un nodo, $\text{root}(G)$, y para cada árbol del conjunto se agrega una rama desde $\text{root}(G)$ hacia la raíz del árbol, que por ser raíz no tenía ramas entrantes. Entonces tenemos una estructura conexa con un nodo que no tiene ramas entrantes y cualquier otro nodo tiene exactamente una rama entrante. Por lo tanto esa estructura es un árbol.

Con el Algoritmo 7 se resume todo el procedimiento. La entrada es una secuencia de n elementos que son m -uplas de símbolos de \mathcal{A} . En la aplicación para dos pistas, m es 2. En la eliminación de ruido m es 1. En esa descripción contexts genera las dos cadenas de formación de contextos, y stree es una función que construye el árbol de sufijos de una secuencia.

Algoritmo 7 optimal – tree (\mathbf{x}) : $[\mathcal{A}^m]^n$

- 1: **Salida:** $T_{opt} : \mathcal{T}_{opt}(\epsilon, \epsilon)$
 - 2: $(\mathbf{y}_\ell, \mathbf{y}_r) \leftarrow \text{contexts}(\mathbf{x})$
 - 3: $(T_{y_\ell}, T_{y_r}) \leftarrow (\text{stree}(\mathbf{y}_\ell), \text{stree}(\mathbf{y}_r))$
 - 4: **Para** $i = 1$ to n
 - 5: $g(\mathbf{s}_{ri}) = (\mathbf{s}_{\ell i}, \mathbf{s}_{ri})$, donde $(\mathbf{s}_{\ell i}, \mathbf{s}_{ri}) \in \mathcal{S} \subset \mathbf{L}(T_{y_\ell}) \times \mathbf{L}(T_{y_r})$
 - 6: $G \leftarrow \text{graph}(\epsilon)$
 - 7: $T_{opt} \leftarrow \text{prune}(G)$
 - 8: **devuelve** T_{opt}
-

La correctitud del algoritmo es enunciada en el siguiente teorema.

Teorema 6.5 (Correctitud). *Las hojas del árbol T_{opt} generado por el Algoritmo 7 constituyen un conjunto óptimo de contextos bi-direccionales para \mathbf{x} .*

Demostración. Teniendo en cuenta que, según el Teorema 6.4, $\text{graph}(\epsilon)$ construye correctamente el CBDCG, incluyendo split , como se vio al describir el podado, obtenemos un árbol con peso mínimo. Las hojas de este árbol constituyen el conjunto óptimo de contextos. \square

6.4. Complejidad del Algoritmo

En la Sección 5.2.2 se obtuvo una cota de la cantidad de contextos representativos. Vamos a ver que esa también es una cota de la complejidad del Algoritmo 7.

Teorema 6.6 (Complejidad). *La complejidad en el peor caso del algoritmo de obtención de un árbol óptimo correspondiente a una secuencia de longitud n es*

$$O(n \cdot \min \{h(T_{y_\ell}), h(T_{y_r})\}).$$

Demostración. La línea 2 obtiene las cadenas de formación de contextos recorriendo x con tiempo constante por símbolo. Para la línea 3, hay varias construcciones de árboles de sufijos en tiempo lineal [10], [11], [12], [13].

El subgrafo $g(s_r)$ que corresponde a cada hoja de T_{y_r} tiene un único vértice y se construye en tiempo constante. Entonces, el bucle de la línea 4 es $O(n)$.

La línea 7, el podado, también implica tiempo lineal. Devuelve un árbol que a lo sumo tiene n hojas por lo que su cantidad de vértices es $O(n)$. El valor `split` ya había sido calculado al construir el CBDCG lo cual permite que: (a) procesar cada nodo del árbol requiera tiempo $O(1)$, y (b) se procesen sólo los nodos del árbol devuelto sin recorrer los que son descartados. Por lo tanto se concluye que el tiempo de este algoritmo es $O(n)$.

Por lo tanto, la complejidad del algoritmo propuesto está dominada por la construcción del CBDCG, en la línea 6.

El CBDCG es construido mediante una recorrida en post-orden de T_{y_r} en el Algoritmo 4. Este algoritmo se invoca una vez por cada vértice y todo lo que se hace en él es $O(1)$, excepto la llamada a `union`. Entonces cada nodo s_r de T_{y_r} es procesado por el Algoritmo 5. Una característica clave de este algoritmo es que, como el `lca` puede calcularse en tiempo constante (Apéndice A), la ejecución es $O(1)$ excepto por las llamadas recursivas. Por lo tanto, la complejidad de la construcción del CBDCG es igual a la cantidad de llamadas a la función `union`. Para un contexto s_r dado, las llamadas recursivas generan todos los vértices del grafo cuyas componentes izquierdas forman un contexto representativo con s_r . Por lo tanto para cada s_r todas esas componentes izquierdas son diferentes. Entonces la complejidad de procesar s_r es igual a la cantidad de vértices en $R_L(s_r)$. Entonces el presente teorema queda demostrado por el Teorema 5.20. \square

El peor caso se da cuando la altura del árbol es $n-1$, produciendo un algoritmo cuya complejidad es $O(n^2)$

6.5. Datos experimentales

6.5.1. Consideraciones prácticas

La implementación del árbol de sufijos se hizo mediante árboles binarios con semántica *primer-hijo siguiente-hermano*. Se usa la versión *suffix-links* y *hojas abiertas* [12]. El nodo representa un contexto, el cual se representa mediante la

posición de inicio en la cadena y la longitud (o con las posiciones de inicio y fin). Si el nodo es una hoja alcanza con la posición de inicio porque se sabe que termina al final de la cadena. En el caso de las hojas tampoco hacen falta el puntero al primer hijo ni el suffix-link. Teniendo en cuenta esto se mantienen separadas las estructuras en un arreglo para los nodos internos, `interns`, y otro para las hojas, `leaves`. Se sabe que hay n hojas y a lo sumo n nodos internos. Los punteros, primer-hijo, siguiente-hermano, suffix-link, son posiciones en los arreglos. Para saber si un puntero es a un nodo que está en `interns` o en `leaves` se utiliza uno de los bits de la dirección. Vamos a suponer que una posición se puede representar con 4 bytes, lo que permite procesar cadenas de hasta 2^{31} símbolos (en lugar de 2^{32} ya que se necesita el bit mencionado; más adelante se explica que esta cota se va a limitar a 2^{30}). Por lo tanto para un nodo interno hacen falta 20 bytes correspondientes a 5 posiciones: inicio y fin, primer-hijo, siguiente-hermano y suffix-link. Para las hojas hacen falta 8 bytes correspondientes a 2 posiciones: inicio y siguiente-hermano. Entonces el árbol de sufijos se puede implementar con $28n$ bytes.

Una vez terminada la construcción cada árbol se adapta para su función específica. El árbol derecho se mantiene por niveles para recorrerlo desde abajo hacia arriba. Esta estructura requiere menos espacio. El árbol izquierdo se modifica para soportar consultas de los `lca`, lo cual se describe en el Apéndice A. Se necesita mantener arreglos cuyo tamaño, que llamamos p , es el doble de la cantidad de nodos. Entonces, como p puede llegar a ser $4n$, si se limita el largo de las cadenas a 2^{30} también las posiciones de estos arreglos se pueden representar con 4 bytes. Además de las estructuras específicas para calcular el `lca`, se mantienen otros dos arreglos. En `lengths` se mantiene la longitud del contexto representado por cada nodo. Como esa longitud es a lo sumo n y la cantidad de nodos está acotada superiormente por $2n$ el espacio ocupado por `lengths` es $8n$. En `lexord_leaf` se especifica cuál es la posición relativa en orden lexicográfico de los contextos de cada una de las n hojas en el conjunto de los contextos de los hasta $2n$ nodos. Por lo tanto requiere $4n$ bytes.

El grafo también se construye con un arreglo para los bi-sufijos, `bileaf` y otro para los bi-ramificables, `bibranch`. Pero como no se sabe hasta terminar de construirlo la cantidad de bi-ramificables este último arreglo se debe poder redimensionar. Cada bi-sufijo se representa con una posición por lo que `bileaf` se implementa con $4n$ bytes. Para los bi-ramificables se mantienen 3 listas mediante nodos enlazados, una para hijos izquierdos, otra para hijos derechos y otra para el conteo de ocurrencias de los símbolos. La longitud de cada una de estas listas puede llegar a ser $|\mathcal{A}|$. Los nodos ocupan 5 bytes por lo que cada lista puede ocupar $5|\mathcal{A}|$. Además se mantiene la longitud de cada una de esas listas y otro byte con *banderas*.

En esta descripción se asumió que las posiciones se representan con 4 bytes

y los símbolos del alfabeto con 1. Para una implementación más ajustada, en el Cuadro 6.1, se muestra el espacio en cantidad de bits de las diferentes estructuras. La longitud de algunos arreglos es p que es $2(|\text{nodos}| - 1)$ (y por lo tanto $(p+2)/2$ es la cantidad de nodos), donde $|\text{nodos}|$ es la cantidad de nodos del árbol izquierdo. Su máximo valor posible es $4n$ pero su valor real se conoce en el momento de crear los arreglos. Uno de esos arreglos se lo considera en bloques de longitud $l = (1/2)(\log p)$. h es la altura del árbol derecho, y se conoce en el momento de crear el arreglo `depth`.

<code>interns</code>	$5n \lceil \log n \rceil$
<code>leaves</code>	$2n \lceil \log n \rceil$
<code>nodes</code>	$p \lceil \log \frac{p+2}{2} \rceil$
<code>depth</code>	$p \lceil \log h \rceil$
<code>branch</code>	$\frac{p+2}{2} \lceil \log \mathcal{A} \rceil$
<code>firstpos</code>	$\frac{p+2}{2} \lceil \log p \rceil$
<code>length</code>	$\frac{p+2}{2} \lceil \log n \rceil$
<code>lexord_leaf</code>	$n \lceil \log \frac{p+2}{2} \rceil$
<code>rpow2</code>	$(p/l) \log(p/l) \lceil \log p \rceil$
<code>canonics</code>	$2^{l-1} \frac{1}{2} (l-1)(l-2) \lceil \log l \rceil$
<code>b2c</code>	$(p/l) \lceil l-1 \rceil$
<code>bileaf</code>	$n \lceil \log n \rceil$
<code>bibranch</code>	$ \mathcal{B} (1 + 3 \lceil \log \mathcal{A} \rceil + 3(\lceil \log \mathcal{A} \rceil + \lceil \log n \rceil)) \mathcal{A} $

Cuadro 6.1: Espacio en bits de las estructuras

6.5.2. Resultados experimentales

Este algoritmo fue implementado y aplicado en el marco de la aplicación eliminación de ruido al texto “Don Quijote” de Miguel de Cervantes. El alfabeto consistió en 96 símbolos (salto de línea y los correspondientes a los códigos ASCII desde el 32 hasta el 126). Las estadísticas del CBDCG generado se presentan en el Cuadro 6.2. La primera fila da la longitud de las secuencias. Cada secuencia contiene a la anterior. La segunda fila da la altura del árbol derecho. La tercera da la cantidad de contextos bi-ramificables, o sea, la cantidad de vértices internos del grafo. Este número es proporcional a la complejidad del algoritmo. En la cuarta fila se puede ver que $n \cdot h(T_{y_r})$ es una cota superior de la cantidad de bi-ramificables, tal como se propuso en el Teorema 6.6. Al crecer n la relación entre las dos magnitudes no crece, sino que incluso disminuye levemente.

En el Listado 6.1 se muestra el resultado de la corrida de prueba con la secuencia de longitud 640000. En esta implementación no se tuvo en cuenta cuál de

$n (\times 1000)$	10	20	40	80	160	320	640
$h(\mathbf{T}_{y_r})$	14	14	15	18	22	24	28
$ \mathcal{B} $	16489	34074	73513	161879	353018	759116	1627152
$ \mathcal{B} / (n \cdot h(T_{y_r}))$	0.12	0.12	0.12	0.11	0.10	0.10	0.09

Cuadro 6.2: Magnitud de CBDCG's

los dos árboles era el de menor altura sino que de manera predeterminada se usó el árbol derecho para la recorrida y el izquierdo para el cálculo de los lca. En el segmento *Construcción de la estructura del grafo* del listado se muestra para cada nivel del árbol derecho, desde el más profundo hasta la raíz, la cantidad de hojas y nodos internos del árbol y la cantidad de bi-ramificables que tienen a los nodos de ese nivel como componentes derechos. La cantidad de bi-sufijos es igual a la cantidad de hojas. En el segmento *Perfil del conjunto óptimo* se muestra la cantidad de hojas del árbol óptimo que corresponden a cada nivel del árbol derecho, desde la raíz hasta el más profundo.

Esta corrida fue realizada en una máquina con arquitectura x86-64 con 4 CPUs de 3093 MHz, con Caché L1d 32K, L1i 32K, L2 256K, L3 3072K, con 4 GB de memoria. El sistema operativo fue Fedora 25 (Workstation Edition), con núcleo Linux 4.13.16-100.fc25.x86_64.

Listado 6.1: Corrida de prueba

```
Alojando memoria para alfabeto y matrices.
Memoria: 149728; Memoria total: 149728
Hojas: 640000
Alojando memoria para construir árboles.
Memoria: 17920028; Memoria total: 18069756

Construyendo árbol izquierdo.
Tiempo: 0.20; Tiempo total: 0.20
Nodos internos: 308744
Altura: 26
Convirtiendo el árbol para cálculos de LCA.
Tiempo: 0.10; Tiempo total: 0.33
Memoria: 37775918; Memoria total: 55845674

Construyendo árbol derecho.
Tiempo: 0.21; Tiempo total: 0.54
Nodos internos: 309629
Convirtiendo el árbol para construcción del grafo.
Tiempo: 0.16; Tiempo total: 0.70
```


Memoria: 11395896; Memoria total: 67241570
Altura = 28

Eliminando mapeo entre hojas.

Memoria: 2560004; Memoria total: 64681566

Eliminando árbol de sufijos.

Memoria: 17920028; Memoria total: 46761538

Iniciando estructuras para el grafo.

Memoria: 86827200; Memoria total: 133588738

Construcción de la estructura del grafo.

Hojas	Inter.	Biramificables
-------	--------	----------------

2	0	0
5	1	1
18	3	4
31	8	11
42	14	17
66	22	31
94	31	61
151	43	86
213	67	132
283	98	204
467	142	301
657	229	453
1029	337	671
1734	539	1021
3442	897	1668
7247	1745	3023
14865	3659	5934
29288	7508	11850
51522	14683	23177
80400	25627	42915
105144	39977	72257
117526	52188	108175
109262	58356	144258
76199	53336	173032
34014	35118	188553
5999	13139	195258
297	1795	203016
3	66	213464
0	1	237579

Cantidad de biramificables: 1627152
Cantidad de aristas: 10561333
Memoria: 52806665; Memoria total: 186395403
Tiempo: 1.26; Tiempo total: 1.96

Eliminando memoria auxiliar.
Memoria: 5002176; Memoria total: 181393227

Ajustando el tamaño de la estructura de biramificables.
Memoria: 0; Memoria total: 181393227

Estableciendo la posición de las hojas.
Tiempo: 0.01; Tiempo total: 1.97

Eliminando estructura para LCA.
Memoria: 37775918; Memoria total: 143617309
Tiempo: 0.01; Tiempo total: 1.97

Calculando vectores de conteo.
Memoria: 950704; Memoria total: 144568013
.....
Tiempo: 0.69; Tiempo total: 2.66

Cantidad de entradas en los vectores de
conteos de los nodos internos: 3508635
Memoria: 17543175; Memoria total: 162111188

Eliminando memoria auxiliar:
Memoria: 950704; Memoria total: 161160484

Calculando peso de los nodos.
Memoria: 21207472; Memoria total: 182367956
.....
Tiempo: 5.39; Tiempo total: 8.05

Eliminando estructuras auxiliares.
Memoria: 14055088; Memoria total: 168312868

Perfil del conjunto óptimo.

Nivel	Cantidad
0	0
1	14
2	943
3	7888
4	21344
5	27829
6	27639
7	24222

8	18542
9	12401
10	7285
11	3850
12	1928
13	918
14	458
15	278
16	168
17	139
18	107
19	99
20	63
21	48
22	42
23	25
24	20
25	22
26	6
27	4

Tamaño del conjunto óptimo: 156282

Eliminando árbol de construcción.

Memoria: 11395896; Memoria total: 156916972

Podado y eliminación de ruido.

Memoria: 960; Memoria total: 156917932

Tiempo: 0.08; Tiempo total: 8.13

Capítulo 7

Conclusión

En esta tesis se desarrolló un algoritmo para construir de manera eficiente conjuntos óptimos de contextos cuya longitud es no acotada. Se generaliza lo anteriormente presentado para contextos de longitud acotada [5], que tiene aplicación en compresión y eliminación de ruido. El algoritmo construye una estructura, el grafo compacto de contextos bi-direccionales que es una generalización de los árboles compactos de sufijos, definido en el Capítulo 6. Es un grafo acíclico dirigido y enraizado, en el que cada vértice que no es sumidero tiene dos conjuntos de aristas salientes. Cada vértice del grafo corresponde a una clase de contextos, más concretamente a un contexto representativo de la clase. Partir en clases de equivalencia el conjunto de contextos que ocurren y tratar con un contexto representativo de cada clase es lo que hace que el algoritmo sea eficiente. Se pasa de una cantidad de orden cúbico de contextos que ocurren a una cantidad cuadrática de contextos representativos.

Todos los contextos de una misma clase comparten propiedades esenciales para los problemas tratados: vector de conteos y como consecuencia el peso. A su vez, el vector de conteos está determinado por las posiciones en las que ocurre el contexto. Todos los contextos de una clase ocurren en las mismas posiciones. En el Capítulo 2 se reformula la notación habitual definiendo el conjunto de índices de un contexto, esto es, las posiciones de la secuencia en las que el contexto ocurre. Se ve que encontrar un conjunto válido de contextos (de los que los conjuntos óptimos buscados son casos particulares) implica una partición de los índices en conjuntos de índices.

En el Capítulo 4 se demuestra que todos los contextos de una clase de contextos tienen una extensión común más corta, que es un contexto de la clase. Esto no es trivial ya que, a diferencia del caso uni-direccional, las clases no tienen un prefijo común más largo. Esos contextos, extensiones comunes más cortas de sus respectivas clases, son los contextos representativos que se corresponden con los vértices del grafo.

Los contextos representativos (vistos como vértices del grafo) son una generalización de los nodos de los árboles de sufijos compactos [10, 11, 12, 13]. Son un subconjunto del producto cartesiano de los nodos de dos de estos árboles: los correspondientes a las cadenas de formación de contextos definidas en el Capítulo 3. Para poder encontrar los contextos representativos de manera eficiente en el Capítulo 5 se definen dos operaciones entre subárboles compactos de un árbol base, concepto también definido ahí. Una de ellas, `glca`, es una generalización de la obtención de los ancestros comunes más profundos de dos nodos de un árbol (`lca` [17, 18]). Es una generalización porque encuentra los `lca` de todos los nodos de un conjunto de árboles. Se describe un algoritmo que los encuentra de manera eficiente. La otra operación, `union`, es la unión de los nodos de subárboles y el `glca` de ellos.

En el Capítulo 6 se describe la construcción del grafo. Se recorre desde las hojas hacia la raíz (recorrido bottom-up) el árbol de sufijos de una de las cadenas de formación de contextos y para cada nodo del árbol recorrido se encuentra, mediante la operación `glca` en el árbol de sufijos correspondiente a la otra cadena de formación de contextos, los contextos con los que forma un contexto representativo. La operación `union` permite obtener la estructura del grafo. Se demuestra que la complejidad del algoritmo es $O(n \cdot h)$, siendo n la longitud de la secuencia a procesar y h el mínimo de las alturas de los árboles de sufijos de las dos cadenas de formación de contextos. En el peor caso esa altura es n , lo que demuestra que la construcción del grafo es en el peor caso $O(n^2)$.

El algoritmo se implementó y con experimentos se comprobó esta cota.

El posterior podado del grafo permite obtener en los nodos que quedan sin ramas salientes el conjunto óptimo de contextos.

Problemas abiertos

El Teorema 6.6 establece que la complejidad de la construcción del grafo está acotada superiormente por el producto de n y la menor de las alturas de los árboles. En el peor caso la altura de un árbol es n , lo cual conduce a un algoritmo $O(n^2)$. Esto ocurre, por ejemplo, si la secuencia consiste en n veces el mismo símbolo. Sin embargo cabe preguntarse qué tan inusual es este caso. Se ha demostrado que la altura promedio de los árboles de sufijos es $O(\log n)$ para algunos modelos de fuentes [19] (informalmente, cuando hay una débil dependencia entre un símbolo y los que ocurrieron d o más posiciones antes, para un cierto entero d). En particular para la aplicación eliminación de ruido, como los dos árboles no son independientes, se puede estudiar si en algún modelo al menos uno de los dos árboles es casi seguramente de baja altura.

La cota superior $n \times h(T_{y,r})$ de cantidad de contextos representativos encontrada en el Lema 5.19 es para el peor caso, en el que al procesar los nodos de cada

nivel de T_{y_r} es necesario involucrar cada hoja en el cálculo de algún lca . En el otro extremo, en el mejor caso el procesamiento de cada nodo T_{y_r} consistiría en un único cálculo de lca . Esto ocurre si los subárboles de T_{y_ℓ} asociados a los hijos del nodo de T_{y_r} son disjuntos. Entonces la pregunta es qué tan ajustada es la cota $n \times h(T_{y_r})$. En la cuarta fila del Cuadro 6.2 se observa que al crecer n decrece la relación entre la cantidad de contextos representativos y $n \times h(T_{y_r})$. Podría haber modelos de fuentes para los cuales, con esta observación y con la anterior el algoritmo sea $o(n \log n)$ en promedio.

El uso de la estructura definida en esta tesis a los problemas de modelado por contextos presenta inconvenientes para la eliminación de ruido porque en esta aplicación no se conoce la cadena original, limpia. Como consecuencia del grafo se obtiene el conjunto de contextos que minimizan la estimación de la pérdida, y no la pérdida verdadera, por lo que no necesariamente se logra maximizar la eliminación de ruido. Se ha demostrado [20] que una longitud de los contextos mayor a $\log n / (2 \log |\mathcal{A}|)$ compromete la calidad de la estimación, y se ha propuesto en consecuencia limitar a ese valor el largo de los contextos, en concordancia con sugerencias anteriores [21]. Resultados de un comienzo de experimentación con esta aplicación confirman estas conclusiones teóricas. Se pueden hacer experimentos más rigurosos para confirmar o no esta hipótesis. Además, esta limitación teórica en la longitud de los contextos hace pensar en el uso de esta estructura, trunca da de manera adecuada, no con el objetivo de mejorar la calidad de los resultados sino con el de mejorar la eficiencia ya que para este propósito se puede aprovechar la compactibilidad dada por los contextos representativos. Hace falta determinar cómo se puede construir eficientemente el grafo a partir de árboles de sufijos con longitudes truncadas ya que deja de existir una correspondencia uno a uno entre las hojas de ambos árboles.

Bibliografía

- [1] Jorma J. Rissanen. A universal data compression system. *IEEE Transactions on Information Theory*, 29(5):656–664, Sep 1983.
- [2] Marcelo J. Weinberger, Jorma J. Rissanen, and Meir Feder. A universal finite memory source. *IEEE Transactions on Information Theory*, 41(3):643–652, May 1995.
- [3] Frans Willems, Yuri Shtarkov, and Tjalling Tjalkens. The context-tree weighting method: basic properties. *IEEE Transactions on Information Theory*, 41(3):653–664, 1995.
- [4] Álvaro Martín, Gadiel Seroussi, and Marcelo Weinberger. Linear time universal coding and time reversal of tree sources via FSM closure. *IEEE Transactions on Information Theory*, 50(7):1442–1468, 2004.
- [5] Erik Ordentlich, Marcelo J. Weinberger, and Tsachy Weissman. Efficient pruning of bi-directional context trees with applications to universal denoising and compression. *Information Theory Workshop, 2004. IEEE*, 0:94–98, 24-29 Oct. 2004.
- [6] Erik Ordentlich, Marcelo J. Weinberger, and Cheng Chang. On multi-directional context sets. *IEEE Transactions on Information Theory*, 57(10):6827–6836, 2011.
- [7] Erik Ordentlich, Marcelo J. Weinberger, and Tsachy Weissman. Multi-directional context sets with applications to universal denoising and compression. *Information Theory, 2005. ISIT 2005. Proceedings. International Symposium on*, 0:1270–1274, 4-9 Sept. 2005.
- [8] Alfredo Viola. Efficient algorithms for universal denoising. In Lars Arge, Robert Sedgewick, and Raimund Seidel, editors, *Data Structures*, number 08081 in Dagstuhl Seminar Proceedings, Dagstuhl, Germany, 2008.

- [9] Fernando Fernández, Alfredo Viola, and Marcelo J. Weinberger. Efficient algorithms for constructing optimal bi-directional context sets. *Data Compression Conference*, 0:179–188, 2010.
- [10] Peter Weiner. Linear pattern matching algorithms. *FOCS 14th Annual Symposium on Switching and Automata Theory*, 0:1–11, 1973.
- [11] Edward M. McCreight. A space-economical suffix tree construction algorithm. *J. ACM*, 23(2):262–272, 1976.
- [12] Esko Ukkonen. On-line construction of suffix trees. *Algorithmica*, 14(3):249–260, 1995.
- [13] Robert Giegerich and Stefan Kurtz. From Ukkonen to McCreight and Weiner: A unifying view of linear-time suffix tree construction. *Algorithmica*, 19(3):331–353, 1997.
- [14] Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory*. John Wiley and Sons, 2006.
- [15] Ragnar Nohre. *Some Topics in Descriptive Complexity*. Linköping studies in science and technology: Dissertations. Linköping Univ., Department of Electrical Engineering, 1994.
- [16] Frans M. J. Willems, Yuri M. Shtarkov, and Tjalling J. Tjalkens. Context weighting for general finite-context sources. *IEEE Transactions on Information Theory*, 42(5):1514–1520, 1996.
- [17] A. V. Aho, J. E. Hopcroft, and J. D. Ullman. On finding lowest common ancestors in trees. In *STOC '73: Proceedings of the fifth annual ACM symposium on Theory of computing*, pages 253–265, New York, NY, USA, 1973. ACM.
- [18] Michael A. Bender and Martin Farach-Colton. The LCA problem revisited. In *LATIN '00: Proceedings of the 4th Latin American Symposium on Theoretical Informatics*, pages 88–94, London, UK, 2000. Springer-Verlag.
- [19] Wojciech Szpankowski. A generalized suffix tree and its (un)expected asymptotic behaviors. *SIAM J. Computing*, 22:1176–1198, 1996.
- [20] Erik Ordentlich, Krishnamurthy Viswanathan, and Marcelo J. Weinberger. Twice-universal denoising. *IEEE Transactions on Information Theory*, 59(1):526–545, Jan 2013.

- [21] Tsachy Weissman, Erik Ordentlich, Gadiel Seroussi, Sergio Verdu, and Marcelo J. Weinberger. Universal discrete denoising: known channel. *IEEE Transactions on Information Theory*, 51(1):5–28, Jan. 2005.
- [22] Dov Harel and Robert Endre Tarjan. Fast algorithms for finding nearest common ancestors. *SIAM J. Comput.*, 13(2):338–355, 1984.
- [23] Baruch Schieber and Uzi Vishkin. On finding lowest common ancestors: simplification and parallelization. *SIAM J. Comput.*, 17(6):1253–1262, 1988.

Apéndice A

Cálculo del lca

En este apéndice se expone cómo encontrar el lca en tiempo constante.

Definición A.1. El Problema LCA consiste en, dados dos nodos v y w de un árbol enraizado T , encontrar el nodo u de T que es el ancestro común más profundo (lowest common ancestor) de v y w .

Si T es un árbol de contextos resolver este problema permite encontrar el prefijo común más largo de los contextos representados por v y w . El Problema LCA fue planteado en 1973 [17]. La forma en que aquí se resuelve está basada en obras posteriores [18, 22, 23].

Para poder encontrar el lca de dos nodos en tiempo $O(1)$ se realiza un pre-procesamiento del árbol. Suponemos que el árbol tiene n hojas y por lo tanto $O(n)$ nodos. El resultado del preprocesamiento podría consistir en una tabla que mantenga el lca de cada par de nodos. El inconveniente de este enfoque es que tal tabla requeriría espacio y tiempo de construcción $O(n^2)$. Aquí se va a describir cómo reducir esa cota a $O(n)$.

Para el propósito de esta tesis se resuelve un problema extendido que consiste en encontrar el lca de un conjunto de nodos y el símbolo que identifica la rama por la cual cada uno de ellos desciende de ese lca

$$\text{lca_branches} : \text{SET}(\text{node}) \rightarrow \text{node} \times \text{SET}(\mathcal{A}) \quad (\text{A.1})$$

en donde está implícita la estructura del árbol, y las dos listas tienen el mismo tamaño, que es mayor o igual a 2 y menor o igual al tamaño del alfabeto \mathcal{A} .

El Problema LCA se transforma en el Problema RRMQ (restricted range minimum query).

Definición A.2. El Problema RMQ consiste en, dadas las posiciones i y j de un arreglo A de números, encontrar la posición del elemento más pequeño en el sub-arreglo $A[i \dots j]$.

Si el mínimo se presenta en más de una posición aquí elegimos la última de ellas.

El Problema RRMQ es una versión del Problema RMQ en el que el arreglo debe cumplir una restricción por la cual el valor de cada posición debe ser ± 1 el valor de la posición anterior.

Para transformar el problema LCA en el problema RRMQ se hace una recorrida euleriana del árbol (depth first). Cada vez que se visita un nodo, tanto en el sentido padre-hijo como en el sentido hijo-padre, se agrega una entrada en el arreglo cuyo valor es la profundidad del nodo. Dos posiciones consecutivas del arreglo corresponden a los extremos de una rama y por lo tanto cumplen la restricción requerida por el problema. Llamemos `depth` al arreglo. La primera y la última posición corresponderían a la raíz del árbol. La última no hace falta incluirla. En la Figura A.1 los nodos del árbol están numerados de acuerdo a la recorrida euleriana.

Para cada nodo hay una posición en `depth` por cada rama que incide en el nodo. Por cada hoja hay una posición y por cada nodo interno distinto de la raíz hay una por cada hijo y otra más para volver al padre. Para la raíz hay una por cada hijo. El tamaño de `depth` es $2 \cdot (|\text{nodos}| - 1)$. Para ver esto hay que tener en cuenta que la cantidad de ramas es $|\text{nodos}| - 1$ y que cada rama se recorre dos veces, una en cada sentido.

Para encontrar el `lca` de dos nodos del árbol es necesario poder encontrar en `depth` una posición correspondiente a cada uno de ellos. Esto significa que hay que mantener una estructura que, dado un nodo del árbol encuentre una posición que le corresponda, por ejemplo la primera, en `depth`. Llamemos `firstpos` a esta estructura que es un arreglo con una posición por cada nodo. El `lca` buscado tiene en el rango definido por las dos posiciones (incluidas) al menos una posición. Todas las posiciones del rango corresponden a descendientes del nodo buscado, por lo que todas las que tengan el valor mínimo corresponden al `lca`. Hay que encontrar una de ellas, por ejemplo la última, y hay que poder determinar a qué nodo corresponde. O sea hace falta un mapeo de $[1 \dots p] \rightarrow [1 \dots |\text{nodos}|]$, donde p es la cantidad de posiciones de `depth`. Esto se resuelve con un arreglo del mismo tamaño que `depth`, al que llamamos `nodes`. En la figura, el rango de `depth` definido por los nodos 4 y 9 va desde la posición 5 hasta la 15 y en ese rango hay dos posiciones con valor 0 que corresponden como se ve en `nodes` al nodo 1, que es el `lca`.

Falta determinar cuál es la rama por la cual cada nodo desciende del `lca` encontrado. Se mantiene una estructura que asocie cada nodo distinto de la raíz con la rama por la cual desciende de su padre. Es un arreglo que llamamos `branch`. En `nodes` y `depth`, para cada nodo interno hay una posición para cada rama que desciende de él. La posición siguiente indica en `nodes` cuál es el hijo que desciende por esa rama. Con `branch` se obtiene la rama por la cual cada

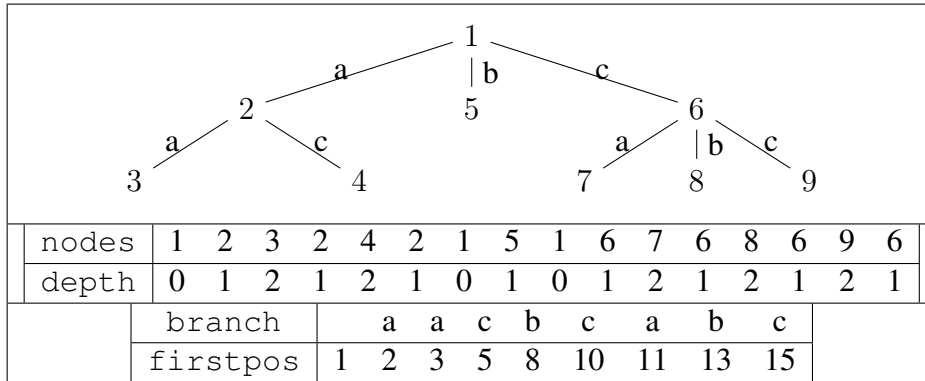


Figura A.1: Transformación del problema LCA en el problema RRMQ

nodo descende de su padre. Si el nodo u es ancestro del nodo v , el lca de ambos es u , y aparece primero en `nodes`. En el rango determinado por ambos nodos las posiciones que en `depth` tienen valor mínimo corresponden a u . La última de estas posiciones (que, tal como se anticipó es la que se obtiene mediante `rrmq`) corresponde al inicio de la rama que es la primera en el camino que conecta ambos nodos. El otro extremo de la rama corresponde al nodo que está en la siguiente posición en `nodes`. Este nodo es el hijo de u que es ancestro de v . Al conocer ese nodo, se obtiene en `branch` la rama buscada. Por ejemplo en el rango definido por los nodos 1 y 8 hay 3 posiciones en `depth` con valor 0. La posición siguiente a la última corresponde en `nodes` al nodo 6, que como se ve en `branch` descende por la rama c .

Para resolver la Función (A.1) sea N la lista de nodos. Supongamos que es un arreglo cuyas entradas son enteros que identifican nodos del árbol. En orden lexicográfico el primero y el último de esos nodos son $a = \min_i \{\text{firstpos}[N[i]]\}$ y $b = \max_i \{\text{firstpos}[N[i]]\}$.

Por lo tanto, recordando el Lema 4.22, el lca de los nodos de N es

$$lca = \text{nodes}[\text{rrmq}(a, b)]$$

y las ramas por las que descenden son

$$B[i] = \text{branch}[\text{nodes}[1 + \text{rrmq}(f, \text{firstpos}[N[i]])]]$$

siendo $f = \text{firstpos}[lca]$.

Por lo tanto hay que hacer una llamada a `rrmq` por cada elemento de la lista de entrada. Como el tamaño del alfabeto es acotado la complejidad de esta operación es $O(1)$ si también lo es el cálculo de `rrmq`.

RRMQ Al arreglo `depth` de tamaño p se lo parte en $m = p/l$ bloques de longitud l . Se verá que conviene que la magnitud de l sea $O(\log p)$.

En tiempo $O(p)$ se obtiene la posición del mínimo de cada bloque. Se construye una tabla que devuelva la posición del mínimo entre dos bloques si la distancia entre ellos es potencia de 2. Es decir, en esa tabla está la posición en `depth` en donde está el mínimo entre el bloque i y los bloques $i + 1, i + 2, i + 4, \dots$. El espacio que ocupa esta tabla es $O(m \cdot \log m)$. Se la puede construir mediante programación dinámica también en tiempo $O(m \cdot \log m)$. La llamamos `rpow2` y tiene una fila por cada potencia de 2 y una columna por cada bloque. Si la longitud de los bloques es $l = O(\log p)$, entonces el espacio ocupado por la tabla y el tiempo necesario para construirla es $O(p)$.

Supongamos que se quiere encontrar la posición del mínimo entre los bloques i y j incluidos. Sea $k = \lfloor \log(j - i) \rfloor$. De `rpow2` se conoce la posición del mínimo entre los bloques i e $i + 2^k$, y la posición del mínimo entre los bloques $j - 2^k$ y j . El menor de los valores en `depth` de las dos posiciones determina cual de los dos es la posición del mínimo en todo el rango.

En el Problema RRMQ se quiere encontrar la posición del mínimo entre las posiciones que en `depth` están en x e y , donde suponemos $x < y$. El bloque, que está entre 0 y $m - 1$, y la posición dentro del bloque, que está entre 0 y $l - 1$, de cada uno de los parámetros es:

$$\begin{aligned} bx &= \lfloor (x - 1)/l \rfloor, & px &= (x - 1) \text{ mód } l, \\ by &= \lfloor (y - 1)/l \rfloor, & py &= (y - 1) \text{ mód } l \end{aligned} \tag{A.2}$$

donde se resta 1 a las posiciones x e y porque están en el rango $[1 \dots p]$.

Cuando dos posiciones pertenezcan al mismo bloque se debe poder resolver el problema en el interior de cada bloque. Para esto para cada bloque se mantiene una tabla de tamaño $\frac{l(l-1)}{2}$, o sea $O(l^2)$. La tabla que corresponde al bloque b es `bloque[b]`, con $0 \leq b \leq m - 1$. En `bloque[b][i][j]` con $0 \leq i < j \leq l - 1$, se encuentra la posición en la que está el mínimo entre i y j en b .

Entonces si x e y pertenecen al mismo bloque el mínimo está en la posición de `depth`

$$1 + bx \cdot l + \text{bloque}[bx][px][py]. \tag{A.3}$$

Si en cambio x e y pertenecen a bloques diferentes el problema se descompone en encontrar las posiciones de tres mínimos y el resultado es la posición en la que está el menor de ellos. Se debe encontrar las posiciones en las que se encuentra el mínimo

- entre x y el fin de su bloque;

- entre el bloque siguiente al de x y el bloque anterior al de y , ambos incluidos;
- entre y y el inicio de su bloque.

El segundo mínimo se obtiene usando la tabla `rpow2`. Se calcula

$$k = \lfloor \log((by - 1) - (bx + 1)) \rfloor,$$

y el mínimo en `depth` está en

$$\begin{aligned} & \text{rpow2}[k][bx + 1] \\ \text{o en} & \text{rpow2}[k][by - 1 - 2^k]. \end{aligned} \tag{A.4}$$

El primero y el tercero son mínimos interiores a los bloques y corresponden respectivamente a las posiciones en `depth`

$$\begin{aligned} & 1 + bx \cdot l + \text{bloque}[bx][px][l - 1] \\ & 1 + by \cdot l + \text{bloque}[by][0][py]. \end{aligned} \tag{A.5}$$

Por lo tanto el mínimo está o bien en la posición dada en la expresión (A.3) si x e y están en el mismo bloque, o bien, si están en bloques distintos, en la posición en la que está el mínimo entre las posiciones dadas por las ecuaciones (A.4) y (A.5).

Como sólo se realiza una cantidad constante de comparaciones y búsquedas en tablas el Problema RRMQ, y por lo tanto también el Problema LCA, se resuelve en tiempo $O(1)$.

Tablas canónicas El inconveniente se presenta al intentar cumplir con el objetivo de espacio requerido y tiempo de preprocesamiento $O(n)$ porque en principio debería haber m tablas, una para cada bloque.

Pero si al valor de cada posición de un bloque se resta el valor de la posición inicial se obtiene otro bloque que empieza con 0. El mínimo entre dos posiciones en este bloque está en la misma posición que el mínimo entre las mismas posiciones en el bloque original. Debido a la restricción que debe cumplir el arreglo original por la cual el valor de cada posición es ± 1 el valor de la posición anterior, sólo hay 2^{l-1} posibles bloques que comienzan con 0, llamémosles *canónicos*. Entonces, en lugar de $m = p/l$ sólo hacen falta 2^{l-1} tablas de tamaño $O(l^2)$. Si $l = \frac{\log p}{2}$, el espacio que ocupan las $\frac{1}{2}\sqrt{p}$ tablas es $O(\sqrt{p} \log^2 p)$, que es $o(p)$. Al conjunto de estas tablas le llamamos `canonics`.

Por ejemplo si $l = 3$ hay 4 tablas canónicas: $[0, -1, -2]$, $[0, -1, 0]$, $[0, 1, 0]$ y $[0, 1, 2]$. La tabla `depth` de la Figura A.1 queda partida en bloques de rangos $[3i + 1 \dots 3(i + 1)]$, con $i \geq 0$. La tabla canónica $[0, 1, 0]$ representaría en `depth` al bloque $[0, 1, 0]$ que está entre las posiciones 7 y 9, pero también a los dos bloques $[1, 2, 1]$ uno de los cuales está entre las posiciones 4 y 6, y el otro entre las posiciones 10 y 12.

Como se ve, estas tablas no dependen de la secuencia a procesar, sino sólo del largo de los bloques. Esto implica que si se pretende procesar varias secuencias de igual longitud estas tablas sólo deben calcularse una vez.

Además es necesario un mapeo que a cada bloque le asigne su tabla canónica. Se resuelve con un arreglo de tamaño m al que llamamos `b2c`. El mapeo de cada bloque a su tabla canónica se hace en tiempo $O(l)$, por lo que todo el mapeo se completa en tiempo $O(p)$.

En las ecuaciones (A.3) y (A.5) se debe sustituir las ocurrencias de `bloque[.]` por `canonics[b2c[.]]`.

Con la construcción de `canonics` y `b2c` queda terminado el preprocesamiento que se hace en tiempo y espacio $O(p)$. Recordando que para esta tesis $p = O(n)$, siendo n el tamaño de la secuencia procesada (igual a la cantidad de hojas del árbol) se cumple el objetivo planteado.