

AED II

The FurfiOS Corporation

Junio 2021

Índice general

1. Especificación	3
1.1. Introducción	3
1.2. Comportamiento Automático	8
2. Diseño	9
2.1. Análisis de Complejidad	9
2.2. Diseño jerárquico de TADs	14
3. Estructuras de Datos	19
3.1. Conjuntos y Diccionarios	19
3.2. Árbol Binario de Búsqueda (ABB)	20
3.3. Árboles AVL	22
3.4. Radix Searching	25
3.4.1. Árboles de búsqueda digital	25
3.4.2. Tries	26
3.5. Hashing	30
3.5.1. Concatenación	32
3.5.2. Direccionamiento abierto	32
3.6. Colas de Prioridad	34
4. Sorting	37
4.1. Heap Sort	38
4.2. Merge Sort	38
4.3. Quick Sort	39
4.4. Bin Sorting	40
4.5. Árboles de decisión	41
5. Dividir y Conquistar	42
5.1. Método de sustitución	43
5.2. Árbol de Recurrencia	44
5.3. Método Maestro	45
6. Memoria Secundaria	47
6.1. Ordenamiento Externo basado en Sort-Merge	47
6.1.1. Fusión Múltiple Equilibrada	48
6.1.2. Selección por sustitución	49
6.1.3. Fusión Polifásica	51
6.2. Búsqueda Externa	51
6.2.1. Acceso Secuencial Indexado	52
6.2.2. Árboles B	52
6.2.3. Hashing Extensible	56
7. Alternativas a los Árboles Balanceados	58
7.1. Skip Lists	58
7.2. Splay Trees	59

Este apunte fue hecho en base a las clases teóricas de los profesores Dr. Esteban Feuerstein y Dr. Emmanuel Iarusi del Primer Cuatrimestre 2020, complementado con bibliografía:

- Capítulos 2, 3, 4, 5.6, 5.7 del Brassard [5].
- Capítulos 4, 11 del Cormen [3].
- Capítulos 4.7, 5.3, 8 del Aho [11].
- Capítulos 13, 15, 17, 18 del Sedgewick [4].
- Capítulo 6.2.4 del Knuth [10].
- Capítulo 4.3 del Tarjan [7]

También se utilizó parte de los apuntes de AED3 (para lo de modelo uniforme y modelo logarítmico) y material extraoficial.

Capítulo 1

Especificación

1.1. Introducción

Típicamente, nos enfrentamos a la siguiente situación. Tenemos un **problema** que se nos es presentado en una manera difusa, vaga, y se pretende que resolver a través de la computadora, de forma eficiente. Lo problemático es transitar este recorrido que va desde la definición informal del problema a su resolución computacional. Ese camino tiene una serie de pasos que, típicamente, las llamamos: **especificación**, **diseño** e **implementación**. Esta serie de pasos son los que vamos a recorrer a lo largo de esta materia.

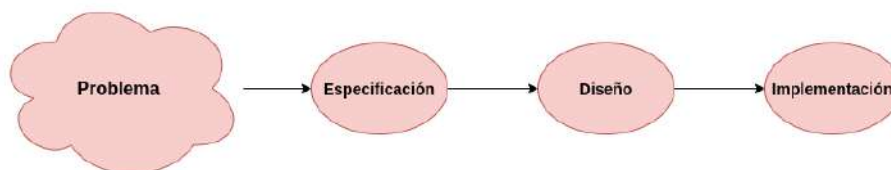


Figura 1.1: Etapas en la resolución de un problema

En esta materia vamos a estudiar **algoritmos** y estructuras de datos. Recordemos qué es un algoritmo:

Definición 1.1.1. *Un algoritmo es un procedimiento para resolver un problema, descrito por una secuencia ordenada y finita de pasos bien determinados que nos llevan de un estado inicial a uno final en un tiempo finito.*

Un algoritmo siempre debe brindar una respuesta, siendo posible que la respuesta sea que no hay respuesta. La descripción debe ser clara y precisa, sin dejar lugar a la utilización de la intuición o la creatividad. Por ejemplo, una receta es un algoritmo si no dice "sal a gusto".

Tenemos que entender que el algoritmo es una herramienta para resolver problemas (como opuesto a recitarlos), de manera tal de poder enfrentar un problema más complejo y poder combinar y modificar distintos algoritmos. Pero, ¿qué es resolver un problema? Se trata de, dada su descripción, proponer un algoritmo que pueda resolver cualquier instancia del mismo.

Los lenguajes naturales tienen una *semántica difusa* (polisemia, ambigüedad). Por lo tanto, al momento de resolver un problema, es necesario tener una descripción precisa del mismo, una que pueda interpretarse sin lugar a dudas, descrita en algún lenguaje formal de especificación rigurosos, con una *semántica precisa*, sin caer en la *sobre-especificación*. Es sencillo realizar una descripción precisa y libre de ambigüedades si damos todos los detalles acerca de la representación del problema. Esta observación nos lleva a la noción de *ocultamiento de la información*. Este salto entre un lenguaje formal y uno natural se conoce como *brecha semántica*.

Metodológicamente, especificar un problema nos resulta útil para empezar a entender la idea de solución algorítmica, por lo que cuando trabajamos con problemas complejos, la solución algorítmica siempre comienza en la especificación. A través de los mecanismos de especificación que vamos a estudiar, vamos a ir introduciendo algunos conceptos importantes como la **abstracción**, la **modularización**, que son conceptos útiles para encontrar soluciones algorítmicas. Notemos que cuando estamos especificando un problema nos interesa describir "el qué" y no "el cómo", por lo que no tiene sentido preguntarnos en esta etapa si la especificación es o no **eficiente**. Una de las claves de la abstracción es poder olvidarse de la representación física de los datos, y tomar a las operaciones como la propia definición de la estructura.

La herramienta principal que utilizaremos para especificar problemas se conoce como Tipos Abstractos de Datos (**TADs**), que son modelos matemáticos que se construyen con el fin de exponer los aspectos relevantes del problema a resolver. Un **tipo abstracto** es un conjunto de valores, los cuales desconocemos su "forma", asociados a operaciones que nos permiten obtener información sobre los valores. Este lenguaje axiomático permite estudiar otras formas de demostración muy útiles, como la *inducción estructural*, y cuenta con la ventaja de que no requiere de tipos primitivos que deban definirse por fuera del mismo. Además, nos resultan de utilidad al utilizar la abstracción, el encapsulamiento, y se tanto flexibles como generales. Veamos un ejemplo:

Descripción : *El dueño de un restaurant quiere asegurarse de que los pedidos sean atendidos con prolijidad. Los mozos llevan los pedidos hasta la cocina donde los colocan. Cuando el cocinero se libera, saca el primer pedido y prepara el plato indicado. El dueño quiere saber cuál es el próximo plato a preparar, cuántos pedidos atiende el cocinero cada día y cuál fue el día con menos pedidos.*

Podemos notar que hay ciertos elementos en esta descripción del problema que realmente no son necesarios para el modelado del mismo. Por ejemplo, en la descripción del problema se nos habla de un *dueño*, pero este no realiza ninguna operación, por lo que no nos interesa modelarlo. Ahora nos toca realizar un proceso de **abstracción** para identificar aquellos elementos que nos interesa modelar. En particular, este problema podemos identificar el uso de platos, días y un restaurant. ¿Qué podemos decir sobre estos elementos?

- Los días pasan, por lo que vamos a querer contarlos. Para modelar este comportamiento nos alcanza con renombrar el TAD DÍA como NAT: **TAD DÍA ES NAT**.
- Solo nos interesa poder diferenciar los platos entre sí. **TAD PLATO ES STRING**.

¿Por qué tenemos que renombrar estos TADs y no utilizar STRING o NAT directamente? Lo que queríamos hacer era utilizar un lenguaje de especificación que nos permita abstraer ciertos conceptos, por lo que usar STRING directamente no sería correcto. Para el restaurant no tiene sentido un STRING, pero los PLATOS sí. De esta manera, nos alejamos de la manera en la que representamos los datos en una computadora y nos acercamos a la definición del problema.

Para especificar un TAD, debemos definir la *signatura*, que nos dice qué operaciones ofrece el TAD, y los axiomas, que determinan el comportamiento de la operaciones. Pasemos ahora a definir la **signatura** del RESTAURANT, en busca de modelar el comportamiento del restaurant definido de forma informal:

TAD RESTAURANT

géneros restaurant

operaciones

cant_platos_pendientes : restaurant \rightarrow nat

próximo_pedido : restaurant $r \rightarrow$ plato {cantidad_platos_pendientes(r)>0}

preparar_plato : restaurant $r \rightarrow$ restaurant {cant_platos_pendientes(r)>0}

tomar_pedido : restaurant \times plato \rightarrow restaurant

nuevo_día : restaurant \rightarrow restaurant
 día_actual : restaurant \rightarrow día
 platos_por_día : restaurant $r \times$ día $d \rightarrow$ nat $\{d \leq \text{día_actual}(r)\}$
 inaugurar : restaurant \rightarrow día

Fin TAD

Luego, la signatura nos define qué operaciones tiene cada tipo, los parámetros y qué tipos devuelven. Notemos que las operaciones de los TADs son **funciones totales**, o sea, son funciones que tienen que estar definidas para todos los valores del dominio. Por eso, en casos como preparar_plato(), tuvimos que restringir el dominio¹.

Con esto definimos la **sintáctica** del TAD, y lo que nos falta es darle **semántica** o comportamiento a estas operaciones. Para ello, utilizaremos **axiomas** del TAD. Los axiomas son *fórmulas bien formadas*, es decir, un predicado aplicado a *términos*². En muchos casos, los axiomas no será otra cosa que ecuaciones como la siguiente:

$$\text{día_actual}(\text{nuevo_día}(r)) \equiv \text{día_actual}(r) + 1$$

De esta manera, podemos darle un significado a la operación *día_actual* cuando tenemos una entrada de la forma *nuevo_día(r)* (en este caso significa *día_actual(r) + 1*). Estas ecuaciones permiten aplicar operaciones del TAD a términos.

Comenzamos con los axiomas que corresponden al TAD RESTAURANT:

TAD RESTAURANT

axiomas $\forall r$: restaurant, $\forall p$: plato

día_actual(inaugurar()) \equiv 0
 día_actual(nuevo_día(r)) \equiv día_actual(r)+1
 día_actual(tomar_pedido(r,p)) \equiv día_actual(r)
 ...
 cant_platos_pendientes(inaugurar()) \equiv 0
 cant_platos_pendientes(tomar_pedido(r,p)) \equiv cant_platos_pendientes(r) + 1
 cant_platos_pendientes(preparar_plato(r)) \equiv cant_platos_pendientes(r) - 1
 ...
 próximo_pedido(r) \equiv ult(secuencia_de_pedidos(r))
 secuencia_de_pedidos(inaugurar()) \equiv <>
 secuencia_de_pedidos(tomar_pedido(r,p)) \equiv p • secuencia_de_pedidos(r)
 secuencia_de_pedidos(preparar_plato(r)) \equiv com(secuencia_de_pedidos(r))

Notemos que estamos utilizando en esta parte de la axiomatización de RESTAURANT algunas operaciones que corresponden a otro tipo de datos, por ejemplo, al TAD SECUENCIA(α) (ver apunte de TADs

¹No sería correcto utilizar una convención de devolver -1 para indicar que no existe resultado.

²Un **término** son las variables, las constantes y lo que resulta de aplicar las operaciones que no tienen parámetros o bien lo que resulta de aplicar operaciones que sí tienen parámetros sobre otros términos

básicos)³. ...

```
platos_por_día(d,inaugurar())           ≡ 0
platos_por_día(d,tomar_pedido(r,p))     ≡ platos_por_día(d,r)
platos_por_día(d, preparar_plato(r))    ≡ if día_actual(r) = d then
                                         platos_por_día(d,r) + 1
                                         else
                                         platos_por_día(d,r)
                                         fi
platos_por_día(d, nuevo_día(r))         ≡ if día_actual(r)+1 = d then
                                         0
                                         else
                                         platos_por_día(d,r)
                                         fi
```

...

```
(∀d': día) 0 ≤ d' ≤ dia_actual(r) ⇒L
  platos_por_día(r, días_menos_pedidos(r)) ≤ platos_por_día(r, d')
```

Fin TAD

Al momento de axiomatizar, una de las primeras cosas que tenemos que tener en cuenta es que las operaciones que estamos especificando son funciones, así que deberíamos evitar cualquier tipo de inconsistencias. En particular, tenemos que tener en cuenta que no se cuenta con *pattern matching*, dicho de otro modo, todos los axiomas valen a la vez, no debemos especificar sobre los casos restringidos (ya que están fuera del dominio), no debemos *sobre-especificar* (tener varias formas de saber cuál es el resultado para unos valores dados) ni tampoco *sub-especificar* (no decir qué valores toma la función para ciertos valores). Notemos que, en cierto sentido, siempre que utilicemos una restricción, estaremos sub-especificando. Sin embargo, es un uso lícito, ya que no siempre podemos caracterizar los resultados de una operación para todo su dominio.

Además, se espera mantener el encapsulamiento entre los distintos TADs, por lo que si trabajamos con instancias de un TAD en las operaciones de otro, manipularemos esas instancias a través de sus observadores (y no sus generadores). Por ejemplo:

TAD RESTAURANT

...

```
platos_de_cierto_precio(r,c,x) ≡ if ∅?(c) then
                                ∅
                                else
                                if precio(DameUno(c),r) = x then
                                    Ag(DameUno(c), platos_de_cierto_precio(r, SinUno(c), x))
                                else
                                    platos_de_cierto_precio(r, SinUno(c,x))
                                fi
                                fi
```

Es preferible a

```
platos_de_cierto_precio(r, ∅, x) ≡ ∅
```

³ *com* devuelve el comienzo de la secuencia (sin el último elemento), *ult* el último elemento

```

platos_de_cierto_precio(r, Ag(p, c), c) ≡ if precio(p, r) = x then
    Ag(platos_de_cierto_precio(r, c, x))
else
    platos_de_cierto_precio(r, c, x)
fi

```

Fin TAD

Ahora, veamos más en detalle las distintas secciones que tiene un TAD:

- **Parámetros formales**
- **Géneros:** un género es el nombre que recibe el conjunto de valores del tipo. Hay una sutil diferencia entre el nombre del TAD y el género. Si se quiere, pensar en el monoide conmutativo $(\mathbb{N}, +)$ (TAD) y en el conjunto de los números naturales (género).
- **Usa:** la sección **Usa** de una TAD hace referencia a las operaciones y géneros de otros TADs que utiliza el TAD que queremos definir. Estas operaciones y géneros tienen que estar exportados en el otro TAD. Si se usa todo lo que aparece mencionado, podemos obviar esta cláusula.
- **Exporta:** indica las operaciones y géneros que se deja a disposición de los usuarios del tipo. Esta cláusula tiene un valor por omisión: los géneros, los observadores básicos y los generadores.
- **Igualdad Observacional:** La igualdad observacional es un predicado entre instancias del tipo que nos dice cuándo son iguales, desde el punto de vista de su comportamiento (semántica) y no desde el punto de vista de la sintáctica ($Ag(1, Ag(2, \emptyset))$ es observacionalmente igual a $Ag(2, Ag(1, \emptyset))$). Para definir la igualdad observacional utilizaremos ciertas funciones para ver detalles particulares de las instancias. Esas funciones se conocen como **observadores básicos**. La $=_{obs}$ es un predicado del metalenguaje, y permite agrupar a las distintas instancias en una misma clase de equivalencia. De esta manera, todas las instancias que sean equivalentes de acuerdo con la igualdad observacional deben mantener la *congruencia* al aplicar cualquier función. Recordemos que una función f es congruente con respecto a una relación de equivalencia " \sim " si y solo si: $(\forall x, y)(x \sim y \iff f(x) \sim f(y))$.
- **Generadores:** son aquellas operaciones que permiten generar o construir instancias del TAD. Es necesario que el conjunto de generadores esté bien definido, es decir, que entre todos los generadores podemos construir cualquier instancia posible del TAD. Un problema menor, no tan grave, es que una instancia del TAD pueda ser construida de más de una manera. Es importante notar que al aplicar un generador sobre una instancia de un TAD **no se está modificando** la instancia que se recibe como parámetro, sino que se genera una **nueva instancia** basada en la anterior. Recordemos que estamos trabajando con *funciones*, por lo que no existe la noción de *estado*, y que el paradigma funcional trabaja bajo el concepto de *transparencia referencial*, es decir, que los resultados de las funciones solo dependen de sus argumentos.
- **Observadores básicos:** son aquellas operaciones que nos permiten diferenciar o distinguir instancias del TAD en clases de equivalencia. En general, se axiomatizan en base a todos los generadores.
- **Extiende**
- **Otras operaciones:** son el resto de las operaciones que se necesitan declarar en un TAD. No debería ocurrir que una función que aparezca en esta sección devuelva valores que rompan con la *congruencia* del TAD (si esto ocurre, habría que repensar los observadores). En general, se axiomatizan en base a los observadores, aunque es posible que en algunos casos sea conveniente axiomatizarlos en base a los generadores.
- **Axiomas:** son las reglas que describen el comportamiento desde el punto de vista semántico de los objetos del TAD.

En general, es preferible que el conjunto de generadores y el de los observadores sean *minimales*. Si estos no lo fuesen, se corre el riesgo de producir inconsistencias y, además, la redundancia atenta contra

la claridad. Típicamente, los observadores deben ser minimales, pero los generadores pueden no serlo (siempre que esto facilite la axiomatización).

1.2. Comportamiento Automático

La idea del comportamiento automático es no modelar operaciones para casos que se dan de forma implícita o automática. Por ejemplo, si cada vez que se da cierta condición A se produce el efecto B a través de una acción C que se da de *forma automática*, seguramente no haga falta hacer alusión a la acción C de ninguna forma para modelar correctamente el objeto de estudio. Veamos un ejemplo.

Descripción del Problema: Se quiere especificar el comportamiento de una fábrica de empanadas que está totalmente automatizada. A medida que se encuentran listas, las empanadas van saliendo de una máquina una a una y son depositadas en una caja para empanadas. En la caja caben 12 empanadas y cuando esta se llena, es automáticamente despachada y reemplazada por una caja vacía. Se quiere saber cuántas cajas de empanadas se despacharon en total, y cuántas empanadas hay en la caja que está actualmente abierta.

El enunciado nos dice que cuando la caja actual se llena es *automáticamente* despachada y reemplazada por una caja vacía. Por lo tanto, no debemos definir una operación de *despacharCaja*, ya que estaríamos permitiendo que esto no ocurra. El mayor impacto que tiene el comportamiento automático sobre la especificación de un problema es en los axiomas, ya que este comportamiento debe quedar plenamente descrito en los mismos. En general, cuando alguna parte del comportamiento del TAD debe ser automática, no deberíamos especificar una acción *manual* ára este comportamiento, no deberíamos permitir que existan instancias del TAD en las que el comportamiento debería haberse aplicado y no lo hizo, y tampoco deberíamos restringir acciones que requieran que suceda el comportamiento, ya que este es automático.

Capítulo 2

Diseño

2.1. Análisis de Complejidad

Hasta ahora, estuvimos enfocándonos en modelar correctamente el problema con algoritmos que terminen y que sean correctos. De esta manera, describimos, mediante un lenguaje formal, *qué* es lo que necesitamos resolver. Vamos a orientarnos, ahora, hacia el *cómo* vamos a resolver el problema, es decir al *diseño* de la solución. Para *diseñar* o elegir la forma adecuada de implementar un TAD, tenemos que tener en cuenta el concepto de la **eficiencia**, de manera tal que esta medida de eficiencia nos permita elegir entre distintos algoritmos que permiten resolver el mismo problema o, en términos de TADs, distintas formas de implementar un TAD.

Cuando hablamos de *eficiencia*, nos referimos a cuántos recursos requiere el algoritmo para ejecutar. En general, nos va a interesar poder medir el **tiempo de ejecución** de un algoritmo, pero también hay otros recursos que habitualmente se usan, como pueden ser el uso de **memoria**, cantidad de procesadores, utilización de la red de comunicaciones, etc. Entonces, debemos formalizar esta idea de consumir más o menos recursos, para así poder diferenciar las distintas implementaciones de forma objetiva y no "mirando a ojo". Otros criterios sobre la eficiencia de los algoritmos pueden ser la claridad de la solución y la facilidad para su codificación.

A partir de estos criterios, nos podríamos preguntar si es posible optimizar todos estos criterios al mismo tiempo. En general, esto no va a ser posible, ya que el uso de estos recursos suele entrar en conflicto, es decir, vamos a tener un *trade-off* entre la utilización de los distintos recursos. Por lo tanto, en cada caso particular, vamos a tener que balancear o decir qué criterio nos importa más, y así poder elegir con qué algoritmo nos quedamos. El recurso que más no va a interesar es el tiempo de ejecución y, en segundo lugar, el espacio utilizado. Esto se debe a que el tiempo no puede ser recuperarse, mientras que el espacio en la memoria sí.

Una primer estrategia es la **empírica** (*a posteriori*), que consiste en programar las distintas soluciones y probarlas con la ayuda de una computadora, para un conjunto arbitrario de instancias. Para medir la complejidad algorítmica vamos a utilizar una medida **teórica**, que nos permita *estimar* lo que tardaría la ejecución de un algoritmo, sin tener que ejecutarlo. Este enfoque teórico permite realizar el análisis de la complejidad *a priori*, aún antes de escribir el programa, el cual vale para todas las instancias, es independiente del lenguaje de programación y de la máquina en la que se ejecuta. Una última alternativa para analizar la eficiencia de un algoritmo es una estrategia *híbrida* entre las dos anteriores, donde la forma de la función que describe la eficiencia del algoritmo es determinada de forma teórica (si es un polinomio, si es exponencial, logarítmica, etc.), y luego a todos los parámetros numéricos requeridos (por ejemplo, los coeficientes del polinomio) son determinados de forma empírica, usualmente a través de algún tipo de regresión¹.

Lo primero que necesitamos para independizarnos de una computadora en particular, es tener un

¹Notemos que hacer este tipo de extrapolaciones, basándonos únicamente en un pequeño conjunto de instancias, puede llevar a predicciones imprecisas e incluso erróneas.

modelo de cómputo. Es decir, vamos a inventar una máquina teórica, ideal, cuyas características son consensuadas, y vamos a asociar la complejidad a esa máquina teórica. Además, vamos a dar una medida de complejidad en función del **tamaño** de las instancias, y no de instancias particulares. Sin embargo, es posible que no todas las instancias de un mismo tamaño se comporten de la misma manera, por lo que estudiaremos las distintas complejidades para cada **tipo** de entrada. Por último, en busca de poder brindar información para todas las instancias de un problema, vamos a enfocarnos en un análisis **asintótico** de la complejidad.

Modelo de cómputo

Estamos buscando una medida que sea general, válida para distintas implementaciones del algoritmo e independiente de la máquina en la que se ejecuta. Con este objetivo en mente, vamos a inventar una máquina teórica que vamos a usar como "banco de pruebas" para la ejecución (teórica) del algoritmo. Esta máquina ideal nos va a permitir definir los conceptos de tiempo de ejecución y espacio de memoria utilizado. Vamos a entender al tiempo de ejecución como la cantidad de pasos o instrucciones que se ejecutan en la máquina teórica para resolver una instancia del problema. Por otro lado, vamos a utilizar como medida del consumo de memoria en función de la cantidad de posiciones de memoria utilizados en la ejecución del algoritmo en la máquina teórica.

Para definir una unidad de tiempo en la máquina teórica, vamos a utilizar el concepto de **operaciones elementales**. Las operaciones elementales son aquellas que "tardan" una cantidad constante de unidades de tiempo en ejecutarse en el modelo de máquina que estamos definiendo. En general, vamos a tener un conjunto reducido de operaciones elementales, cuyo tiempo de ejecución no depende del tamaño de la entrada, y reglas para calcular cuánto tardan aquellas operaciones que no son elementales.

En teoría, estas operaciones podrían no ser elementales, ya que el tiempo de ejecución aumenta con la longitud de los operandos. Sin embargo, en general es posible hacer esta asunción, porque los operandos envueltos en las instancias que esperamos encontrar son de una medida *razonable*. Sin embargo, en nuestro modelo teórico el algoritmo debería poder ser aplicado a todas las instancias posibles, y como ninguna máquina real podría ejecutar en tiempo constante sumas y multiplicaciones de instancias indefinidamente grandes, en algunos contextos es necesario hacer un análisis más fino para este tipo de operaciones.

Típicamente vamos a tener dos modelos. El **modelo uniforme** asume que cada operación básica tiene un tiempo de ejecución constante, y resulta apropiado cuando los operandos de las operaciones se restringen a una palabra de la máquina. En cambio, en el **modelo logarítmico**, el tiempo de ejecución de cada operación es una función (que dependerá del algoritmo utilizado para resolverla) del tamaño (medido en **bits**) de los operandos. Luego, una suma resulta proporcional al tamaño del mayor operando, mientras que una multiplicación resulta proporcional al producto de los tamaños (se puede mejorar). Este modelo es apropiado cuando los operandos de las operaciones intermedias pueden crecer arbitrariamente.

En general, vamos a trabajar bajo el modelo uniforme y consideraremos como operaciones elementales a las operaciones aritmético-lógica básicas (suma, división, multiplicación, AND, OR), las comparaciones lógicas, las transferencias de control y las asignaciones a variables de tipos básico.

Vamos a utilizar una función de complejidad $t(I)$, que mide la cantidad de operaciones elementales que se ejecutan cuando un algoritmo se ejecuta para una instancia I . Luego, para poder calcular esta función de complejidad, es necesario tener en cuenta algunas consideraciones generales. Vamos a considerar que el tiempo de una OE es de una **unidad**.

En general, para el análisis del tiempo de ejecución de un algoritmo vamos a ver cuánto cuestan las operaciones individuales, para luego combinar estos costos según la estructura de control en la que se encuentren las operaciones. Entonces, para poder analizar la complejidad de un algoritmo, debemos saber cómo lidiar con **estructuras de control** y **ecuaciones de recurrencia**.

Estructuras de Control:

- Secuencias de instrucciones: Si P_1 y P_2 dos fragmentos sucesivos de un algoritmo, podemos calcular el costo total del algoritmo como $t(A) = t(P_1) + t(P_2)$. Notemos que P_1 y P_2 podrían ser instrucciones simples o complicados sub-algoritmos.

- CASE (switch en C), donde la expresión es evaluada una vez y comparada con cada una de las constantes asociadas al case, y si resultan iguales, las instrucciones subsiguientes al case son ejecutadas, hasta encontrar un break; si no la expresión resulta distinta a todas las constantes, se ejecuta el default statement.

Vamos a considerar que el tiempo de ejecución de este tipo de sentencias se puede expresar como: $T = T(E) + T(S_i)$, es decir, lo que tardamos en evaluar la expresión E (incluyendo el tiempo de comparación con las distintas constantes), más lo que tardamos en evaluar alguno de los códigos que haya que ejecutar (S_1, S_2, \dots, oS_n), dependiendo del valor de E. Cuando pensamos en el análisis del peor caso, tenemos que reemplazar $T(S_i)$ por el $\max\{T(S_1), T(S_2), \dots, T(S_n)\}$.

- Cuando consideramos el tiempo de ejecución de la sentencias **if E then S₁ else S₂**, lo podemos expresar como $T = T(E) + \max\{T(S_1), T(S_2)\}$ (para el análisis del peor caso). Además, cuando consideramos el tiempo de ejecución de sentencias del tipo **while E do S**, lo podemos expresar como

$$T = T(E) + \prod_{i=1}^{\text{iteraciones}} T(S_i) + T(E_i).$$

Notemos que tanto $T(E)$ como $T(S)$ pueden variar en cada iteración, y por tanto habrá que tenerlo en cuenta para el cálculo final.

- Por último, el tiempo de ejecución de una llamada a un procedimiento $F(P_1, P_2, \dots, P_n)$ es

$$T = 1 + T(P_1) + T(P_2) + \dots + T(P_n) + T(F),$$

ya que tenemos un costo 1 por la llamada a F, más el tiempo de evaluación de los parámetros, más el tiempo que tarde en ejecutarse F. No vamos a contabilizar la copia de los parámetros al stack, sino que asumimos que estos se pasan por referencia. Si este no fuera el caso, habría que tener en cuenta la copia de los parámetros.

En resumen, podemos obtener el tiempo de ejecución de una secuencia consecutiva de instrucciones calculando sumando los tiempos de ejecución de cada una de las instrucciones y estructuras de control. Notemos que el tiempo de ejecución de las llamadas a procedimientos recursivos va a dar lugar a ecuaciones en recurrencia ($T(n) = n + T(n - 1)$), que veremos posteriormente (ver Divide & Conquer).

Tamaño de la entrada

Habíamos dicho que íbamos a definir la complejidad de un algoritmo en función del tamaño de entrada, pero ¿por qué es importante? Lo que queremos es una complejidad relativa al tamaño de entrada porque los valores que podemos medir para instancias particulares dependen del tipo de máquina, del contexto, lo cual es poco generalizable. Queremos una medida que sea más general y que tenga valor para muchas instancias. Ese es el motivo por el que definir los recursos en función del tamaño de entrada.

El tamaño de entrada formalmente se corresponde con el número de bits necesarios para representar a esa instancia en una computadora, usando algún esquema de codificación. Sin embargo, normalmente vamos a ser menos formales que esto, y usaremos a la palabra *tamaño* para significar cualquier entero que de alguna manera mida el número de componentes de una instancia. Por ejemplo, cuando hablamos del problema de ordenamiento, medimos el tamaño de la instancia según el número de elementos a ordenar. A veces, cuando hablamos de problemas que involucran enteros, podemos dar la eficiencia del algoritmo en términos del *valor* de la instancia, en lugar de su tamaño (notar que la cantidad de bits que ocupa un entero de valor n es $\lceil \log n \rceil$).

Para ello, vamos a definir a la función $T(n)$, que nos habla de la complejidad temporal de un algoritmo en función del tamaño n de una entrada (y no de la entrada propiamente dicha), y a la función $S(n)$ para la complejidad espacial. Ahora, esta definición tiene el problema de que no todas las entradas del mismo tamaño consumen el mismo tiempo (o espacio), por lo que es necesario afinar un poco esta medida. Esto da a lugar a tres medidas particulares para un mismo algoritmo: el análisis del caso **peor**, del caso **mejor** y del caso **promedio**.

Definición 2.1.1. Sea $t(I)$ el tiempo de ejecución de un algoritmo sobre una instancia I. Definimos el tiempo de ejecución del peor caso, del mejor caso y del caso promedio para instancias de un tamaño n como:

$$\begin{aligned}
T_{peor}(n) &= \max_{|I|=n} t(I) \\
T_{mejor}(n) &= \min_{|I|=n} t(I) \\
T_{promedio}(n) &= \sum_{|I|=n} P(I) \cdot t(I)
\end{aligned}$$

donde $|I|$ indica que la instancia I es de tamaño n .

Intuitivamente, $T_{peor}(n)$ es el tiempo de ejecución del algoritmo sobre la instancia que implica mayor tiempo de ejecución entre las entradas de tamaño n , mientras que $T_{mejor}(n)$ nos habla del tiempo de ejecución del algoritmo sobre la instancia que implica menor tiempo de ejecución. Por último, el análisis del caso promedio es una medida muy utilizada, y nos habla del tiempo de ejecución esperable sobre instancias "típicas" (ponderado por la probabilidad de ocurrencia de cada instancia).

El problema con el análisis promedio es que requiere de conocer la distribución estadística de las entradas, lo cual en muchos casos no es realista. Además, tomar un promedio puede ser engañoso, por ejemplo, podemos tener algoritmos para los cuales tenemos muchas entradas que requieren 1 paso, y muchas otras que requieren 100 pasos, dándonos un promedio de 50, a pesar de que ninguna entrada tarde 50 pasos. Por lo tanto, la medida que más vamos a utilizar es la del caso peor, al ser la medida más robusta, nos da garantías y nos permite hacer afirmaciones certeras sobre el comportamiento de los algoritmos, aunque estas sean pesimistas. Es decir, vamos a medir la eficiencia de los algoritmos por lo peor que les puede suceder cuando tengan que ejecutarse sobre una instancia de un determinado tamaño.

Comportamiento Asintótico

En general, cuando hacemos un análisis de la complejidad de un algoritmo, nos va a interesar poder determinar su **comportamiento asintótico**. Esta idea se basa en el **principio de invarianza**², que nos dice que dado un algoritmo y dos implementaciones M_1 y M_2 que tienen un tiempo de ejecución $T_1(n)$ y $T_2(n)$, siendo n el tamaño de la entrada, existen $c \in \mathbb{R}^+$ y $n_0 \in \mathbb{N}$ tales que:

$$T_1(n) \leq c \cdot T_2(n), \quad \forall n \geq n_0$$

Por lo tanto, no nos va a interesar tanto conocer exactamente la cantidad de operaciones de un algoritmo, sino más bien el **orden de magnitud** del tiempo de ejecución. La idea, entonces, va a ser buscar aquella función (logarítmica, lineal, cuadrática, exponencial, etc.) que exprese el comportamiento del algoritmo, para entradas suficientemente grandes. Para ello, se han propuesto distintas medidas del comportamiento asintótico: \mathcal{O} (cota superior), Ω (cota inferior) y Θ (orden exacto de la función).

Esta notación se dice *asintótica* porque trabaja sobre el comportamiento de funciones en el límite, es decir, para valores suficientemente grandes de sus parámetros. En consecuencia, los argumentos basados en la notación asintótica podrían fallar en cuanto a su valor práctico cuando trabajamos con entradas con valores del "mundo real". En cualquier caso, suele ser de utilidad al momento de comparar algoritmos, y nos facilita algunas cuentas. Notemos que la utilización de las cotas asintóticas para comparar funciones de tiempo de ejecución se basa en la hipótesis de que son suficientes para decidir el mejor algoritmo, prescindiendo de las constantes de proporcionalidad. Sin embargo, esta hipótesis puede no ser cierta cuando el tamaño de la entrada es pequeño, o cuando las constantes involucradas son demasiado grandes, en cuyo caso mantendremos el coeficiente del término de mayor peso (pensar en si conviene usar un algoritmo cúbico en segundos o un algoritmo cuadrático en días).

La notación \mathcal{O} sirve para representar el límite o cota superior del tiempo de ejecución de un algoritmo. Es decir, la notación $f \in \mathcal{O}(g)$ expresa que la función f **no crece** más rápido que alguna función proporcional a g , y decimos que g es cota superior de f . Luego, si sabemos que un algoritmo tiene un tiempo de ejecución $T_{peor} \in \mathcal{O}(g)$, podemos asegurar que para todas las entradas de tamaño suficientemente grande, el tiempo $T_{peor}(n)$ va a ser proporcional a g . Formalmente, dada una función $f: \mathbb{N} \rightarrow \mathbb{R}^{\geq 0}$ arbitraria que va desde los números naturales a los reales no negativos, por ejemplo siendo

²Este principio no es algo que se pueda probar, sino que se sustenta en base a la observación.

$f(n) = n^2$, siendo n un representante del tamaño de la instancia del algoritmo y $f(n)$ la cantidad de recursos utilizados por el algoritmo para procesar esa instancia. Luego, decimos que $f \in \mathcal{O}(g)$ si y solo si

$$\exists n_0, c > 0 \text{ tal que } n \geq n_0 \Rightarrow f(n) \leq c \cdot g(n)$$

Por conveniencia, permitimos el uso incorrecto de notación para decir que $f(n)$ está en el orden de $g(n)$ incluso si $f(n)$ es negativo o indefinido para una cantidad finita de valores de n , y solo vamos a requerir que esté bien definida cuando $n \geq n_0$.

El *umbral* n_0 suele ser útil para simplificar argumentos, pero nunca es necesario cuando consideramos funciones *estrictamente* positivas. Sean $f, t : \mathbb{N} \rightarrow \mathbb{R}^+$ dos funciones que van de los naturales a los reales positivos. La **regla del umbral** nos dice que $f(n) \in \mathcal{O}(g(n))$ si y solo si existe una constante real positiva c tal que $f(n) \leq c \cdot g(n)$ para *cada* número natural n . La idea es que si tenemos que $f(n) \leq c \cdot g(n)$ para todo $n \geq n_0$, podemos probar que también vale para $n_0 = 0$. Sea b la razón máxima entre $f(n)$ y $g(n)$, con $0 \leq n < n_0$ (sabemos que el máximo existe pues es un conjunto finito). Entonces, $f(n) \leq b \cdot g(n)$ para los $0 \leq n < n_0$. Luego, si tomamos $k = \max(b, c)$, habremos obtenido una constante real positiva tal que $f(n) \leq k \cdot g(n)$ para todo n natural. Esta regla se puede generalizar para decir que si $f(n) \in \mathcal{O}(g(n))$ y $g(n)$ es positiva para todo $n \geq n_0$, entonces ese mismo n_0 puede usarse como *umbral* para la notación \mathcal{O} .

La herramienta más poderosa y versátil para probar que cierta función está en el orden de otra se conoce como la **regla del límite**, que nos dice que, dadas dos funciones arbitrarias f, g ambas $\mathbb{N} \rightarrow \mathbb{R}^{\geq 0}$:

1. Si $\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} \in \mathbb{R}^+$, entonces $f(n) \in \mathcal{O}(g(n))$ y $g(n) \in \mathcal{O}(f(n))$. La vuelta no vale, ya que es posible que el límite no exista.
2. Si $\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = 0$, entonces $f(n) \in \mathcal{O}(g(n))$ pero $g(n) \notin \mathcal{O}(f(n))$.
3. Si $\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = +\infty$, entonces $f(n) \notin \mathcal{O}(g(n))$ pero $g(n) \in \mathcal{O}(f(n))$.

Una herramienta útil para probar que una función está en el orden de otra es la **regla del máximo**. Esta nos dice que $\mathcal{O}(f(n) + g(n)) = \mathcal{O}(\max(f(n), g(n)))$.

Consideremos el problema de ordenamiento. Ya sabíamos que existen algoritmos como Insertion Sort o Selection Sort que pueden ordenar un arreglo de n elementos en $\mathcal{O}(n^2)$. Sin embargo, existen otros algoritmos más sofisticados como *Heap Sort* que tienen costo $\mathcal{O}(n \cdot \log n)$. Está claro que $n \log n \in \mathcal{O}(n^2)$. Por lo tanto, sería correcto decir que *heapsort* está en $\mathcal{O}(n^2)$, o incluso $\mathcal{O}(n^3)$. Esto se debe a que la notación \mathcal{O} está diseñada solamente para dar cotas superiores acerca de la cantidad de recursos requeridos. Claramente, necesitamos una notación dual para dar una cota *inferior*. La notación Ω justamente nos permite representar el límite o cota inferior del tiempo de ejecución del algoritmo.

Consideremos nuevamente dos funciones $f, g : \mathbb{N} \rightarrow \mathbb{R}^{\geq 0}$ que vayan de los naturales a los reales no negativos. Luego, la notación $f \in \Omega(g)$ expresa que la función f está acotada inferiormente por alguna función múltiplo positivo de g para valores de n lo suficientemente grandes. Formalmente, decimos que $f \in \Omega(g)$ si y solo si

$$\exists n_0, k > 0 \text{ tal que } n \geq n_0 \Rightarrow f(n) \geq k \cdot g(n).$$

Podemos ver que $f(n) \in \Omega(g(n))$ si y solo si $g(n) \in \mathcal{O}(f(n))$ (esta propiedad se conoce como **regla de dualidad**).

A pesar de la gran similitud entre la notación \mathcal{O} y la notación Ω , hay un aspecto en el que la dualidad falla. Recordemos que normalmente vamos a querer estudiar el tiempo de ejecución en el *peor caso*. Por lo tanto, cuando decimos que $T_{\text{peor}}(n) \in \mathcal{O}(g(n))$ para el peor caso, estamos diciendo que **para toda** instancia de tamaño n , existe una constante real positiva k tal que $T_{\text{peor}}(n) \leq k \cdot g(n)$. En cambio, cuando decimos que $T_{\text{peor}}(n) \in \Omega(g(n))$, estamos diciendo que existe **al menos una** instancia de tamaño n para la cual realmente tiene costo $t(I) \geq k \cdot g(n)$ (para todo $n \geq n_0$). Por lo tanto, podrían existir infinitas instancias de tamaño n que se podrían resolver en un menor tiempo.

En general, vamos a utilizar la notación Ω para dar cotas inferiores en los tiempos de ejecución de

algoritmos. Sin embargo, también es posible dar cotas inferiores a la dificultad propia de resolver cierto problema. Por ejemplo, vamos a ver que para **cualquier** algoritmo que sea capaz de ordenar n elementos, basándose en comparaciones de a pares de elementos, se tiene que pertenece a $\Omega(n \cdot \log n)$. Luego, se dice que el *problema* de ordenamiento por comparaciones tiene una complejidad en $\Omega(n \cdot \log n)$. Este resultado es mucho más fuerte que decir que cierto algoritmo particular tiene una cota inferior, e incluso nos dice cuándo un algoritmo de ordenamiento es **óptimo**.

Finalmente, nos queda la notación Θ (orden exacto), que vamos a definir como el conjunto de funciones que crecen asintóticamente de la misma forma, es decir: $\Theta(f) = \mathcal{O}(f) \cap \Omega(f)$. Formalmente:

$$\Theta(g) = \{f \mid \exists n_0, k_1, k_2 > 0 \text{ tal que } n \geq n_0 \Rightarrow k_1 \cdot g(n) \leq f(n) \leq k_2 \cdot g(n)\}$$

Tanto la regla del umbral como la regla del máximo siguen valiendo para la notación Θ . Además, las reglas de límite se pueden reformular de la siguiente forma:

1. Si $\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} \in \mathbb{R}^+$, entonces $f(n) \in \Theta(g(n))$.
2. Si $\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = 0$, entonces $f(n) \in \mathcal{O}(g(n))$ pero $f(n) \notin \Theta(g(n))$.
3. Si $\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = +\infty$, entonces $f(n) \in \Omega(g(n))$ pero $f(n) \notin \Theta(g(n))$.

Es posible que cuando analicemos el costo de un algoritmo, este dependa simultáneamente de más de un solo parámetro. En estos casos, la noción de "tamaño de entrada" que estuvimos usando hasta ahora pierde un poco de sentido. Por esta razón, la notación asintótica se generaliza para funciones de varias variables. Sea $g : \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{R}^{\geq 0}$ una función que va de los pares de números naturales a los reales no negativos, por ejemplo $g(m, n) = m \cdot \log n$. Sea $f : \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{R}^{\geq 0}$ otra función. Decimos que $f(m, n)$ está en el orden de $g(m, n)$, denotado por $f(m, n) \in \mathcal{O}(g(m, n))$, si $t(m, n)$ está acotada superiormente por un múltiplo positivo de $g(m, n)$ cuando tanto m como n son suficientemente grandes. Formalmente, $\mathcal{O}(g(m, n))$ se define como

$$\{t : \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{R}^{\geq 0} \mid (\exists c \in \mathbb{R}^+) (\forall^\infty m, n \in \mathbb{N}) [t(m, n) \leq c \cdot g(m, n)]\},$$

donde \forall^∞ quiere decir que la propiedad vale para todo natural, exceptuando por una cantidad finita de excepciones (vale para enteros lo suficientemente grandes). Para el caso general de n parámetros, la definición es similar.

Fórmula de Stirling : Nos sirve para aproximar factoriales grandes, y por tanto reescribir aquellas complejidades que dependan de una suma de logaritmos. La aproximación se expresa como

$$\ln(n!) \approx n \cdot \ln(n) - n$$

Luego, si un algoritmo $A \in \Theta(\log n!) \Rightarrow A \in \Theta(n \cdot \log n)$.

2.2. Diseño jerárquico de TADs

Ahora nos toca ver cómo pasamos de la *especificación* ("el qué") al *diseño* ("el cómo") del algoritmo. Estuvimos viendo que para una misma especificación, podemos diseñar distintos algoritmos. Para saber qué algoritmo nos conviene utilizar, deberemos considerar aspectos como la eficiencia en el tiempo de ejecución y el espacio utilizado, de acuerdo con el *contexto de uso*, es decir, en qué orden llegarán los datos, cómo se los consultará cuáles serán las operaciones más frecuentemente usadas, cuál debe ser su complejidad en el peor caso, etc. Además, se deben tomar en cuenta el uso de las buenas prácticas, como el encapsulamiento y el ocultamiento de la información. Al diseñar un algoritmo bajo una cierta especificación, estamos pasando de un paradigma *funcional* a un paradigma *imperativo*, y queremos que ese pase resulte lo más ordenado, metódico posible. Notemos que cuanto más abstracto sea el TAD que vamos a implementar, más opciones de diseño tendremos disponibles.

Tenemos como objetivo obtener un diseño jerárquico y modular. Para ello, el método que veremos tiene las nociones de los distintos niveles en la jerarquía. Cada uno de los niveles tendrá asociado un módulo de abstracción. Es decir, habrá distintos tipos abstractos de datos que deberemos diseñar, y a cada uno de ellos le corresponderá un módulo de abstracción. Por ejemplo, consideremos al TAD CONJUNTO y estos dos posibles diseños:

- Un arreglo redimensionable, con una Inserción (sin repetidos y en orden) en $\mathcal{O}(n)$ y una Búsqueda en $\mathcal{O}(\log n)$.
- Una secuencia, con una Inserción en $\mathcal{O}(1)$ y una Búsqueda en $\mathcal{O}(n)$.

Podríamos preguntarnos cuál de estos dos diseños nos conviene, pero esto va a depender del *contexto de uso* en el cual vamos a utilizar esta estructura de datos (¿vamos a tener muchas más búsquedas que inserciones?).

A grandes rasgos, el método se compone de la elección del TAD a diseñar, el módulo de abstracción para el TAD elegido, e ir iterando hasta finalizar entre los distintos niveles. En general, el orden en el cual se diseñan los TADs es arbitrario, pero resulta una buena práctica comenzar por los más importantes o de mayor nivel de abstracción, al ser estos quienes imponen los requerimientos de eficiencia para los TADs menos importantes. Decimos que esta modalidad de trabajo es *top-down*. La idea va a ser, por ejemplo, comenzar diseñando un modelo para representar al TAD CONJUNTODeNat, en el cual utilizamos valores del TAD SECUENCIADeNat, que se podrían representar con punteros.

Además, contamos con tipos de datos **primitivos** (`bool`, `nat`, `int`, `real`, `char`, `puntero`) para nuestro lenguaje de diseño los cuales se asume que ya están definidos. También tienen la particularidad de que los parámetros de tipos primitivos son los únicos que siempre se pasan **por valor**, mientras que los tipos no primitivos se pasan por **referencia** (salvo que se diga lo contrario, en cuyo caso se pasa una copia y se asume el costo de realizar dicha copia).

Cada módulo de abstracción está compuesto por dos secciones: la definición de la *interfaz* y la definición de la *representación*. En la interfaz se describe la funcionalidad del módulo y en qué contexto puede ser usado: complejidad temporal de los algoritmos, *aliasing* y efectos secundarios de las operaciones. Recordemos que estamos trabajando bajo el paradigma imperativo, donde los datos son tratados como entidades independientes de las funciones que los utilizan, por lo que dejamos de tener *transparencia referencial*. Además, suele pasar que tengamos distintas variables que hacen referencia a una misma estructura física, a lo que llamamos *aliasing*.

Para poder formalizar un poco esta idea de aliasing entre dos variables, vamos a definir un "metapredicado", *alias*(ϕ), que nos va a servir para describir aquellas variables que comparten memoria en la ejecución del programa. *alias* es un metapredicado con un único parámetro ϕ que es un predicado que involucra un conjunto V de dos o más variables del programa. El significado que le damos es que las variables de V satisfacen ϕ durante la ejecución del resto del programa, al menos hasta que le asignemos otro valor a dichas variables. También puede ocurrir que la variable quede indefinida, por ejemplo, cuando el valor al que hace referencia deja de ser válido (normalmente ocurre cuando se elimina un elemento que está siendo iterado). La idea es que al modificar cualquier variable en V , el resto de las variables se deben actualizar manteniendo el invariante ϕ . Por ejemplo, si tenemos *alias*($s = t$), estamos diciendo que s y t comparten la misma memoria de forma tal que cualquier modificación sobre s afecta a t y viceversa, satisfaciendo $s = t$ mientras que a s o a t no se les asigne otro valor. En general, vamos a usar aliasing para iteradores, punteros y referencias a una misma estructura física.

Cuando estamos pasando del mundo funcional al mundo imperativo, nos van a aparecer algunas cuestiones que tenemos que resolver desde el punto de vista formal. En particular, cuando estamos definiendo la Interfaz, vamos a querer definir funciones de la siguiente forma:

```

PERTENECE(in C : conjuntoDeNat, in n : nat) → res : bool
Pre ≡ {true}
Post ≡ {res =obs n ∈ C}

```

Sin embargo, en esta expresión estaríamos mezclando variables del mundo imperativo con conceptos del mundo funcional (como la igualdad observacional). Para resolver este problema de manera formal,

vamos a introducir un operador funcional que, dado un valor en el mundo imperativo, nos permite vincularlo el valor correspondiente del mundo funcional. Es decir, introduciremos una notación definida como una función que va del género imperativo G_I al género teórico G_T :

$$\hat{\bullet} : G_I \rightarrow G_T$$

Luego, reescribimos a la operación *pertenece* como:

```

PERTENECE(in C : conjuntoDeNat, in n : nat) → res : bool
Pre ≡ {true}
Post ≡ {rés =obs n̂ ∈ Ĉ}

```

En la sección de *representación* se elige una forma de representación utilizando otros módulos, y se resuelven las operaciones del módulo en función de su representación. También es necesario definir cuál es la relación entre la estructura de representación y el TAD representado, para finalmente definir, explicar y mostrar cómo son los algoritmos para implementar las operaciones del módulo, utilizando elementos de módulos de menor jerarquía, incluyendo todos los cálculos que justifican las complejidades de las operaciones.

La *estructura de representación* describe los valores sobre los cuales se representará el género que se está implementando. Esta estructura de representación tiene que ser capaz de representar todos los posibles valores del TAD y, además, debe poder describir cómo es que las operaciones del TAD se traducen en operaciones que se realizan sobre la estructura (satisfaciendo las exigencias del contexto de uso). Para ello, vamos a tener que explicar cómo se relaciona la representación con la abstracción, es decir, qué quiere decir qué.

La estructura de representación de las instancias solo será accesible a través de las operaciones que se hayan detallado en la interfaz del módulo de abstracción, permitiendo separar cada nivel en la jerarquía de diseño, ocultando cómo es que realmente se representa cada tipo.

La siguiente parte de un módulo de diseño son los Algoritmos.

Consideremos el siguiente problema. Se nos pide implementar un conjunto de naturales, bajo el siguiente contexto: *los números del 1 al 100 deben manejarse en $\mathcal{O}(1)$, mientras que el resto en $\mathcal{O}(n)$. Además, se debe poder conocer rápidamente la cardinalidad.* Nuestra propuesta va a ser la siguiente: representar al conjunto en tres partes

- Un arreglo de 100 posiciones booleanas (un 1 en la posición i indica que el elemento i pertenece al conjunto). Esto nos va a permitir manejar a los números del 1 al 100 en $\mathcal{O}(1)$.
- Además, vamos a incluir en nuestra representación una secuencia, donde se almacenarán todos los elementos que pertenezcan al conjunto y que sean mayores que 100.
- Y por último, vamos a incluir un *nat* para la cardinalidad.

Esto lo vamos a notar de la siguiente manera:

conjunto_semi_rápido se representa con estr

donde **estr** es `tupla(rápido: arreglo [1 ... 100] de bool, resto: secu(nat), cardinal: nat)`

Ahora bien, si nos convencimos de que esta es una buena representación, tenemos que preguntarnos si es posible representar a cualquier valor del TAD CONJUNTO y, por otro lado, si cualquier valor que tome esta estructura tiene una contraparte en el TAD.

Está claro que todo conjunto puede ser representado por esta estructura, sin embargo, no todos los valores corresponden a instancias válidas. Por ejemplo, la tupla $([0, \dots, 0], \langle 37, 107, 28 \rangle, 3)$ no representa ninguna instancia del TAD. Por lo tanto, es necesario tener algún predicado que nos permita distinguir entre representaciones válidas de las representaciones inválidas, y así poder establecer una relación en-

tre la representación y el valor representado del TAD. Este predicado se conoce como **invariante de representación**.

Este invariante de representación lo vamos a agregar a las precondiciones, y vamos a tener que cumplir en las postcondiciones de las operaciones del módulo. Esto nos va a obligar, pero también permitir, mantener ciertas garantías sobre las estructuras de representación, lo que nos puede permitir hacer las cosas más eficientemente (podemos usar algoritmos que aprovechen las garantías del invariante).

Volviendo al ejemplo, las condiciones del invariante de representación, descritas de manera informal, deberían ser las siguientes:

1. Que *resto* solo tenga números mayores que 100.
2. Que *resto* no tenga números repetidos.
3. Que *cardinal* tenga la longitud de *resto* más la cantidad de celdas de *rápido* que estén en 1 (*true*).

Veamos ahora la descripción formal del invariante. Formalmente, el invariante de representación es una función booleana que tiene como dominio la versión abstracta del género de representación (el $\hat{\bullet}$), e indica si la instancia es o no válida. Luego, el invariante de representación, descrito formalmente, nos queda:

$$\begin{aligned}
 &Rep : e\hat{str} \rightarrow bool \\
 &(\forall e : e\hat{str}) Rep(e) \equiv true \Leftrightarrow \\
 &\quad (1)(\forall n : nat)(esta?(n, e.resto) \Rightarrow n > 100) \wedge \\
 &\quad (2)(\forall n : nat)(cant_apariciones(n, e.resto) \leq 1) \wedge \\
 &\quad (3)e.cant = long(e.resto) + contar_trues(e.rapido)
 \end{aligned}$$

donde *esta?*, *cant_de_apariciones* y *long* son funciones de secuencias, y *contar_trues* es una función sobre arreglos.

Con esto podemos restringir el conjunto de valores posibles que podría tomar la estructura de representación. Sin embargo, todavía nos falta poder vincular las distintas instancias válidas de la estructura con el valor representado del TAD. Para ello, utilizamos la **función de abstracción**. Esta función de abstracción toma una instancia abstracta ($\hat{\bullet}$) de la estructura de representación y nos devuelve una instancia abstracta del género representado.

La función de abstracción debe ser un homomorfismo con respecto a la signatura del TAD, es decir, para toda operación \bullet del módulo, se tiene que cumplir que:

$$Abs(\bullet(p_1, \dots, p_n)) =_{obs} \hat{\bullet}(Abs(p_1), \dots, Abs(p_n))$$

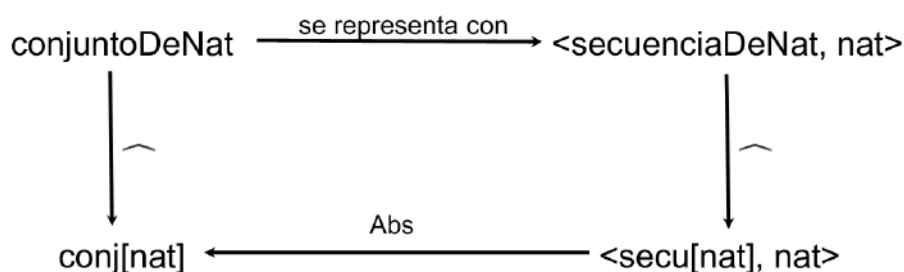


Figura 2.1: Función de Abstracción

Normalmente, tenemos dos formas de describir la función de abstracción. La primera de ellas es en función de sus observadores básicos, dado que estos permiten identificar de manera unívoca a cualquier

instancia. Otra forma de describirla en función de los generadores del TAD. Veamos cómo podemos escribir la función de abstracción del módulo de conjunto_semi_rápido:

Abs: $\hat{\text{estr}}\ e \rightarrow \text{conjunto_semi_rápido}\ \{\text{Rep}(e)\}$

$$\begin{aligned} \text{Abs}(e) =_{\text{obs}} C / \\ (\forall n : \text{nat})(n \in C \Leftrightarrow ((n \geq 100 \wedge r.\text{rápido}[n]) \vee \\ (n > 100 \wedge \text{está?}(n, e.\text{resto}))) \end{aligned}$$

La función de abstracción debe ser total, restringiéndose a las instancias que cumplan con el invariante de representación. Además, no tiene por qué ser inyectiva, es decir, dos estructuras diferentes pueden representar al mismo término de un TAD, pero tiene que ser sobreyectiva sobre las clase de equivalencia del TAD (dadas por la igualdad observacional).

Por último, nos falta diseñar los algoritmos que implementan las distintas operaciones del módulo. Junto con los algoritmos, incluiremos la precondition y la postcondition. Veamos cómo notamos el algoritmo para la operación AGREGAR(...):

AGREGAR(**in/out** C : conjunto_semi_rápido, **in** e : nat)
Pre $\equiv \{\hat{C} =_{\text{obs}} C_0 \wedge \hat{e} \notin C\}$
Post $\equiv \{\hat{C} =_{\text{obs}} \text{Ag}(C_0, \hat{e})\}$

Luego, una posible implementación del algoritmo podría ser la siguiente:

```

1 iAgrega(in/out C : estr, in e : nat)
2 C.cant++
3 if e < 100 then
4   | C.rápido[e] = true
5 else
6   | AgregaAtrás(C.resto, e)

```

Cuando estamos en la etapa de implementación, tenemos la flexibilidad de ponerle un nombre especial para el algoritmo. En este caso, lo llamamos **i**, de implementación, pero podríamos ponerle algún otro nombre. Otra cosa a tener en cuenta es que estamos trabajando directamente sobre la estructura de representación (y no la abstracción).

Capítulo 3

Estructuras de Datos

En los capítulos anteriores estuvimos hablando del camino que hay que recorrer entre la descripción informal hasta llegar a una solución concreta, y vimos algunas herramientas formales para describir y realizar ese proceso. En esta parte nos vamos a dedicar a estudiar soluciones concretas a problemas concretos, que nos van a permitir avanzar en la implementación de problemas fundamentales de la informática.

3.1. Conjuntos y Diccionarios

Vamos a empezar con los TADs CONJUNTOS y DICCIONARIOS, que constituyen uno de los núcleos de la algoritmia: el problema de organizar y buscar la información en un sistema informático. En los conjuntos, lo único que nos interesa es saber si un elemento particular pertenece o no al conjunto, mientras que en los diccionarios, además de saber si una *clave* está definida para ese diccionario, nos interesa poder encontrar su *significado*.

Es fácil ver que si tenemos una implementación de un diccionario, podemos implementar un conjunto simplemente ignorando el significado de las claves. Por otro lado, también podemos implementar un diccionario a partir de un conjunto, por ejemplo, guardando a las claves junto con sus significados en una tupla (clave, significado). En caso de tener como significado a una estructura compleja, podemos guardarnos un puntero a la misma, en lugar de la estructura en sí. Otra opción para el manejo de claves con significados podría ser mantener dos estructuras paralelas en las que la posición i de una tengamos la clave, y en la misma posición de la otra tengamos su significado. Esto nos estaría diciendo que ambos problemas son equivalentes. Luego, vamos a trabajar sobre la representación de este tipo de estructuras, buscando representaciones o implementaciones eficientes de las mismas.

La primer idea que se nos ocurre para representar a los conjuntos es a través de estructuras secuenciales: en cada posición del arreglo, almacenamos la tupla (clave, significado), y un puntero para saber dónde termina el arreglo.

Ahora, pensemos cuáles van a ser los distintos órdenes de complejidad de las operaciones básicas:

- VACÍO: para crear un diccionario vacío, nos basta con crear un arreglo vacío y definir el puntero al último como NIL, por lo que la operación nos cuesta $O(1)$.
- DEFINIR: para definir un nuevo elemento, tenemos que agregar este nuevo elemento al final del arreglo y actualizar el puntero al último elemento, por lo que la operación nos cuesta $O(1)$.
- DEFINIDO?: para saber si una clave está o no definida, tenemos que realizar una búsqueda secuencial por todo el arreglo, por lo que la operación nos lleva $O(n)$.
- SIGNIFICADO: para obtener el significado de una clave, nuevamente tenemos que realizar una búsqueda secuencial, con un costo $O(n)$.

Otra posibilidad es mantener el arreglo ordenado, permitiendo hacer uso de la búsqueda binaria, mejorando la complejidad de DEFINIR? y OBTENER hasta $O(\log n)$. El problema de mantener el arreglo ordenado es que al momento de DEFINIR una nueva clave, debemos buscar la posición que le toca en el arreglo, y correr todos los elementos siguientes, con un costo de $O(n)$. Por lo tanto, este segundo modelo sería mejor en contextos de uso estáticos, en el que se realicen muchas búsquedas y pocas inserciones.

Cualquiera de estas dos opciones tienen operaciones con complejidad en $O(n)$, por lo que nos gustaría tener alguna estructura más eficiente, de manera tal que todas las operaciones se puedan realizar con un costo sub-lineal. Para ello, vamos a representar conjuntos a través de árboles binarios.

En un árbol binario tenemos un nodo inicial, el nodo raíz, del cual pueden salir a lo sumo dos nodos (hijo izquierdo e hijo derecho), y por cada nodo que no sea el raíz, tenemos un nodo padre y podemos tener a lo sumo dos hijos. Pensando un poco sobre el costo de las operaciones de un diccionario, podemos ver que las complejidades serían las siguientes:

- VACÍO: para crear un diccionario vacío, nos basta con crear un árbol vacío, el cual tiene como raíz a un NIL, por lo que la operación nos cuesta $O(1)$.
- DEFINIR: para definir un nuevo elemento, tenemos que definir una regla para agregarlo al árbol. Podemos tomar como regla que se agrega en el primer lugar posible en el que entra el nodo, sin tener que inaugurar un nuevo nivel. Esto lo podemos hacer manteniendo un puntero al último, por lo que la operación nos cuesta $O(1)$.
- DEFINIR?: para saber si una clave está o no definida, tenemos que realizar una búsqueda secuencial por todo el árbol, por lo que la operación nos lleva $O(n)$.
- SIGNIFICADO: para obtener el significado de una clave, nuevamente tenemos que realizar una búsqueda secuencial, con un costo $O(n)$.

Podemos ver que no hemos mejorado en cuanto a la complejidad de las operaciones. Sin embargo, lo que podemos hacer es restringir un poco la forma en la que está organizado el árbol binario, manteniendo un invariante sobre la estructura que nos permita mejorar el costo de las operaciones.

3.2. Árbol Binario de Búsqueda (ABB)

Un grafo es un conjunto de nodos unidos por un conjunto de líneas (no orientadas) o flechas (orientadas), en ambos casos llamados *ejes*. Si todas los vértices están unidos por líneas, decimos que es un grafo no orientado; si todos los vértices están unidos por flechas, decimos que es un grafo orientado (o digrafo). También es posible tener grafos mixtos, donde algunos nodos están unidos por líneas y otros por flechas. Decimos que un grafo es conexo cuando podemos ir de cualquier nodo a cualquier otro nodo siguiendo una secuencia de ejes. Las secuencias de ejes pueden formar caminos y ciclos. Un *árbol* es un grafo no orientado, conexo y acíclico. Si cuenta con un nodo distinguido, llamado *raíz*, decimos que el árbol es *enraizado*. Normalmente, cuando dibujamos un árbol enraizado, ponemos la raíz en la cima, similar a un árbol familiar, con los otros ejes saliendo hacia abajo de la raíz. Podemos describir las relaciones entre los nodos haciendo una analogía del árbol familiar, usando términos como *padre*, *hijos*, hermanos, ancestros, etc. Una *hoja* en un árbol es un nodo sin hijos; los otros nodos se denominan nodos *internos*. Un árbol k -ario es aquel en el que cada nodo puede tener a lo sumo k hijos. En esta sección, vamos a trabajar específicamente sobre árboles binarios, donde cada nodo puede tener 0, 1 o 2 hijos.

Un **Árbol Binario de Búsqueda** (ABB) es un árbol binario que satisface que para todo nodo, los valores de los elementos en su subárbol izquierdo son menores que el valor de ese nodo, y los valores de los elementos de su subárbol derecho son mayores. Luego, podemos definir el invariante de un árbol ABB de la siguiente forma:

1. el valor de todos los elementos del subárbol izquierdo es menor que el valor de la raíz,
2. el valor de todos los elementos del subárbol derecho es mayor que el valor de la raíz,
3. el subárbol izquierdo es un ABB,
4. el subárbol derecho es un ABB.

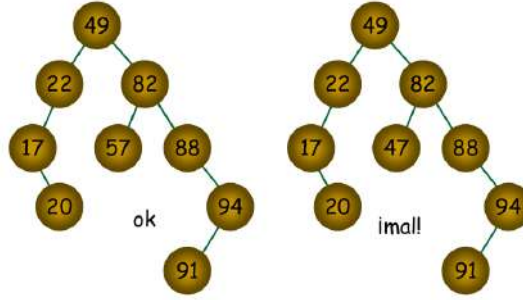


Figura 3.1: Ejemplos de árboles binarios.

Formalmente, asumiendo que tenemos una estructura de representación para los árboles binario $ab(nodo)$, por ejemplo, $\langle I, r, D \rangle$, donde I es el subárbol izquierdo, r la raíz y D el subárbol derecho, podemos plantear el siguiente invariante de representación para los árboles ABB:

$$\begin{aligned}
 ABB : ab(nodo) &\rightarrow bool \\
 (\forall e : ab(nodo)) ABB(e) &\equiv true \Leftrightarrow \\
 &(0) ABB(e) \equiv Nil?(e) \vee \\
 &(1) (\forall c : clave) (esta?(c, Izq(e)) \Rightarrow c < clave(raiz(e))) \wedge \\
 &(2) (\forall c : clave) (esta?(c, Der(e)) \Rightarrow c > clave(raiz(e))) \wedge \\
 &(3) ABB(Izq(e)) \wedge \\
 &(4) ABB(Der(e))
 \end{aligned}$$

Notemos que para poder construir un ABB, es necesario poder dar un orden total¹ a las claves. Ahora, veamos algunos algoritmos básicos para los ABB:

iVacio(in $\langle \rangle$, out $A : ab(nodo)$)

1 return NIL

iDefinir(in $c : clave, s : significado$, in/out $A : ab(nodo)$)

```

1 if  $A = NIL$  then
2    $A \leftarrow ab(NIL, \langle c, s \rangle, NIL)$ 
3 else
4   //  $A = ab(I, \langle r_c, r_s \rangle)$ 
5   if  $c < r_c$  then
6      $A \leftarrow ab(definir(c, s, I), \langle r_c, r_s \rangle, D)$ 
7   else
8      $A \leftarrow b(I, \langle r_c, r_s \rangle, definir(c, s, D))$ 
9 return  $A$ 

```

Podemos ver que el costo de inserción va a depender de la cantidad de recursiones que tenemos que hacer para llegar desde el nodo raíz hasta la hoja donde vamos a colocar el nuevo nodo. Vimos que no es demasiado complicado actualizar un ABB, sin romper su invariante. Sin embargo, si esto se hace de una manera descuidada, el ABB resultante podría quedar *desbalanceado*. Con esto nos referimos a que muchos nodos en el árbol tengan un solo hijo, de manera tal que las ramas se vuelvan largas y finas. Cuando esto ocurre, las operaciones en el árbol dejan de ser eficientes. En el peor caso, todos los nodos en el árbol podrían tener exactamente un hijo, exceptuando la única hoja. Por lo tanto, para la búsqueda de un elemento en el árbol nos quedaría en $O(n)$ en el peor caso, al tener que comparar con todos los nodos del árbol. Por otro lado, si las claves que nos van llegando siguen una distribución uniforme, es

¹Dadas dos claves cualesquiera c_1, c_2 , debemos poder determinar si $c_1 \leq c_2$.

esperable que nos vaya quedando un árbol casi balanceado (casi completo), por lo que nos quedaría un árbol de altura $\approx \log n$, y por tanto la búsqueda nos quedaría $O(\log n)$.

Para el algoritmo de borrado, $\text{Borrar}(u : \text{nodo}, A : \text{ab}(\text{nodo}))$ tenemos tres casos:

1. u es una hoja. En este caso, lo único que tenemos que hacer es buscar al padre de u en A , y luego eliminar la hoja u (reemplazándolo por NIL). Si u no tiene padre, significa que es la raíz, y como es hoja, estamos en el caso de un árbol de un único nodo, por lo que al quitarlo, nos quedaríamos con el árbol vacío.
2. u tiene un solo hijo. En este caso, primero debemos buscar al padre w de u . Si u no tiene padre, significa que u es la raíz, por lo que debemos reemplazar a la raíz de A por la raíz del único subárbol de r_A . Si u tiene padre, tenemos dos casos: o bien u es hijo izquierdo o bien es hijo derecho. Si es hijo izquierdo, significa que u es menor a w , al igual que todos los hijos de u . Entonces, podemos mover a todo el subárbol de u (no importa si es subárbol derecho o izquierdo), y colocarlo donde antes estaba u , ya que sabemos que todos ellos son menores a w . El caso de que u sea hijo derecho es análogo al anterior.
3. u tiene dos hijos. En este caso, lo que vamos a hacer es buscar un reemplazo de u dentro de alguno de los subárboles de u , el cual vamos a mover desde su posición original a la posición de u . Para ello, debemos encontrar algún nodo que cumpla con el invariante que debe cumplir u , es decir, debe ser mayor a todos los nodos del subárbol izquierdo de u y debe ser menor a todos los nodos del subárbol derecho de u . Para cumplir esta condición, tenemos dos posibles candidatos: el nodo de más a la derecha dentro del subárbol izquierdo y el nodo de más a la izquierda dentro del subárbol derecho. Vamos a considerar el primer caso (el otro es similar).

El nodo que está más a la derecha dentro del subárbol izquierdo es mayor a todos los nodos del subárbol izquierdo, además, es menor que u , por lo que también es menor a todos los nodos del subárbol derecho. Por lo tanto, queda claro que lo podríamos mover hacia la posición de u . Para moverlo, lo que vamos a hacer es copiarlo a la posición de u , y luego borrarlo dentro del subárbol izquierdo al que pertenecía. Para ello, tenemos que darnos cuenta que al ser el nodo más a la derecha, no puede tener un hijo derecho, porque sino no sería el de más a la derecha. Luego, estamos en el caso 2, en el que tenemos que borrar un nodo con un solo hijo, que sabíamos resolver.

La propiedad de los árboles ABB nos permite recorrer todas las claves de manera ordenada a través de un algoritmo recursivo, llamado **inorder tree walk**. Se llama así porque primero imprime todas las claves del subárbol izquierdo, luego la clave de la raíz, y finalmente todas las claves del árbol derecho (recursivamente). Similarmente, tenemos el **preorder tree walk** que imprime la raíz antes que cualquier subárbol, y el **postorder tree walk** que imprime la raíz después de ambos subárboles.

INORDER-TREE-WALK(x)

```
1 if  $x \neq \text{NIL}$  then
2   INORDER-TREE-WALK(izq)
3   imprimir  $x$ 
4   INORDER-TREE-WALK(der)
```

Se puede demostrar que si inicialmente x es la raíz del árbol, entonces el costo total del algoritmo nos queda en tiempo $\Theta(n)$ (pensar que solo visitamos a cada nodo una única vez, y que el costo de cada visita es constante).

3.3. Árboles AVL

Queríamos obtener una estructura que nos permita mejorar la complejidad de las operaciones de diccionario para el peor caso, y con los árboles ABB seguimos teniendo el mismo problema de antes, ya que seguimos teniendo algunas operaciones que tienen costo lineal. Uno de los conceptos claves que deberíamos tener en cuenta, a partir del desarrollo de los ABBs, es que el costo de las operaciones depende de la profundidad del árbol, por lo que si pudiéramos controlar esta profundidad, podríamos

mejorar la complejidad de las operaciones. Luego, nuestro objetivo es encontrar una estructura similar a los ABBs que nos asegure que el árbol se mantenga **balanceado**. La idea es que si logramos que las ramas del árbol sean lo más parecidas entre sí o, dicho de otro modo, que el árbol sea lo más completo posible, tendríamos que la rama más larga, que define el caso peor, sería lo más corta posible.

Existen varios métodos para mantener un árbol balanceado, garantizando la eficiencia en las operaciones de búsqueda, inserción y borrado. Entre las técnicas más viejas tenemos el uso de árboles AVL y árboles 2-3; otras estructuras más modernas incluyen a los Red-Black Trees y los Splay Trees. En esta sección vamos a estudiar las propiedades de los árboles AVL para ver cómo consiguen mantener balanceado el árbol.

Cuando hablamos de árbol con **perfectamente balanceado**, a lo que nos referimos es a que todas las ramas tengan casi la misma longitud (dos ramas difieren en a lo sumo una unidad), y que todos los nodos internos tengan ambos hijos. Se puede demostrar que un árbol perfectamente balanceado de n nodos tiene altura $\lfloor \log_2(n) \rfloor + 1$. Además, para este tipo de árboles, se puede probar que la cantidad de hojas n_h es igual a la cantidad de nodos internos n_i más uno: $n_h = n_i + 1$, lo cual nos dice que más de la mitad de los nodos en un árbol perfectamente balanceado son hojas. El problema del *balanceo perfecto* es que no se conocen algoritmos eficientes que permitan mantener este invariante a través de sucesiones de inserciones y borrados, por lo que tendríamos una degradación en la performance.

Lo siguiente que falta hacer es buscar una propiedad que sea más débil que el balanceo perfecto, pero que nos permita tener operaciones eficientes. Se dice que un árbol es **balanceado en altura** si las alturas de los subárboles izquierdo y derecho **de cada nodo** difieren en a lo sumo una unidad. Antes pedíamos que la altura de **todas** las ramas difieran en a lo sumo una unidad, mientras que ahora lo que estamos pidiendo es que la **altura** del subárbol derecho y el izquierdo difieran en a lo sumo una unidad, es decir, que la rama **más larga** de uno y la rama más larga del otro difieran en longitud a lo sumo una unidad. De esta manera, permitimos tener ramas que varíen en más de una unidad en su longitud, por lo que la propiedad es más débil. Con esta idea en mente, definimos el *factor de balanceo (FDB)* de un nodo u como $FDB = h(D_u) - h(I_u)$, siendo D_u el subárbol derecho de u e I_u el subárbol izquierdo. Decimos que un árbol es balanceado en altura si para todo nodo u vale que $-1 \leq FDB(u) \leq 1$. Vamos a asumir que como parte de la estructura de los árboles AVL tenemos precalculado el factor de balanceo para cada nodo.

Lo que queremos ver es que un árbol AVL tiene una altura logarítmica en función de la cantidad de nodos. Para demostrar esto, lo que se hace es ver que para toda altura h , cualquier árbol AVL tiene una cantidad exponencial de nodos. Para ello, nos alcanza con ver que aquel árbol AVL de altura h de menor cantidad de nodos, tiene una cantidad exponencial de nodos en función de la altura. Este tipo de árboles AVL, que tienen el mínimo número de nodos para una altura h , se conocen como **árboles de Fibonacci**, porque se puede probar que la cantidad de nodos de un árbol de Fibonacci $Fibo_h = Fibo_{h-1} + Fibo_{h-2} + 1$, y que $Fibo_h = F_{h+2} - 1$, siendo F_i el i -ésimo término de la sucesión de Fibonacci. Para ver esto, lo que tenemos que saber es que se puede construir un árbol de Fibonacci de altura $i + 2$ a partir de unir dos árboles de Fibonacci: uno de altura i con otro de altura $i + 1$.

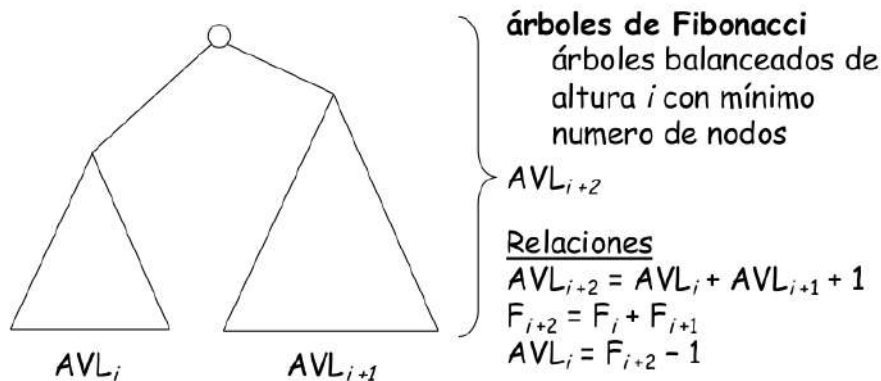


Figura 3.2: Árboles de Fibonacci

A partir de esto, se puede probar que un árbol de Fibonacci con n nodos tiene una altura menor a $1,44 \log(n + 2) - 0,328$, por lo que quedaría demostrado que un AVL de n nodos tiene altura $\Theta(\log n)$, al ser un árbol de Fibonacci lo "más desbalanceado" que puede llegar a ser un árbol AVL. Con todo esto en mente, nos falta ver si es posible realizar las operaciones de diccionarios de forma eficiente, manteniendo esta estructura de los AVLs.

Inserción : Básicamente, la inserción en un árbol AVL consta de tres pasos: primero tenemos que insertar al elemento como lo haríamos en un ABB ($\Theta(\log n)$), recalculamos los FDBs partiendo desde abajo hacia arriba, y luego tenemos que realizar una serie de **rotaciones** que nos permitan recuperar el invariante del AVL, es decir, que $-1 \leq FDB(u) \leq 1$ para todo nodo u .

Para recalcular los FDB, tenemos que darnos cuenta que solo pueden verse alterados son aquellos FDB asociados a los nodos que pertenecen a la rama que va desde la raíz hasta el nodo agregado, por lo que a lo sumo vamos a tener que recorrer $\Theta(\log n)$ nodos, ya que sabemos que la altura de los árboles AVL es logarítmica con respecto a la cantidad de nodos.

Por último, nos queda ver cuántas rotaciones tenemos que hacer para que el árbol vuelva a ser un AVL, y cuánto nos cuesta hacer cada una. Tenemos dos tipos de rotaciones sobre un nodo: rotación a izquierda y rotación a derecha.

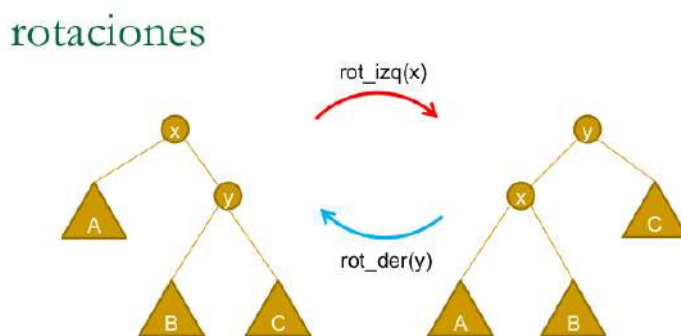


Figura 3.3: Tipos de Rotaciones

Notemos que al momento de aplicar cualquier tipo de rotación, ya sea a izquierda o a derecha, se mantiene el invariante de los árboles ABB, es decir, se cumple que para todo nodo en el árbol, todos los nodos de su subárbol derecho son mayores y todos los nodos de su subárbol izquierdo son menores a este

nodo. Además, podemos ver que una rotación nos cuesta $O(1)$, porque lo único que hacemos es cambiar una cantidad constante de punteros de lugar.

Ahora, lo que nos falta ver es cómo rebalanceamos el árbol que había sido AVL, pero que le agregamos un nuevo elemento. En el peor caso, a lo sumo vamos a tener que hay nodos que tengan $FDB(u) = \pm 2$, con u perteneciente a la rama que va desde la raíz hasta el nodo que queríamos agregar. Para analizar todos los casos posibles, vamos a separar en cuatro casos:

- RR: queremos insertar en el subárbol derecho de un hijo derecho Q . En este caso, nos basta con hacer una rotación a izquierda sobre el padre de Q .
- LR: queremos insertar en el subárbol izquierdo de raíz R de un hijo derecho Q , con padre P . En este caso, tenemos que hacer una rotación doble: primero hacemos una rotación a derecha sobre Q , y por último hacemos una rotación a izquierda sobre P .
- RL: queremos insertar en el subárbol derecho de un hijo izquierdo. Similar al caso LR.
- LL: queremos insertar en el subárbol izquierdo de un hijo izquierdo. Similar al caso RR.

Con todo esto en cuenta, nos queda que el costo del primer y segundo paso es proporcional a la altura del árbol ($\Theta(\log n)$), y el costo del último paso, que consiste en hacer a lo sumo dos rotaciones, nos queda en $O(1)$. Por lo tanto, concluimos que el costo total del algoritmo de inserción en los árboles AVL es $\Theta(\log n)$.

Borrado : El algoritmo de borrado consiste en tres pasos: borrar el nodo como lo haríamos en un ABB, recalculamos los factores de balanceo, y por último tenemos que hacer una rotación simple o doble por cada nodo con $FDB = \pm 2$ ($O(\log n)$ en el peor caso). Por lo tanto, la operación de Borrado también cuesta $\Theta(\log n)$.

3.4. Radix Searching

Existen varios métodos de búsqueda que proceden examinando a las claves de a partes, en lugar de hacer comparaciones entre claves completas en cada paso. Estos métodos, llamado *radix searching methods*. La motivación de estas estructuras es encontrar una implementación de diccionarios que no dependa tanto de la cantidad de claves, pero que mantenga un rendimiento razonable para el peor caso en función de la cantidad de claves. La idea, entonces, es que las operaciones del diccionario terminen dependiendo del *tamaño* de las claves, y no tanto de la *cantidad* de claves. Para ello, nos vamos a enfocar en realizar comparaciones de **partes** de las claves, en lugar de comparar claves completas, por ejemplo, si las claves son enteros, sus partes podrían ser los bits que lo codifican; si las claves son strings, las partes podrían ser sus caracteres.

La principal ventaja de los métodos de radix searching es que proveen una eficiencia razonable en el peor caso sin la complicación de los árboles balanceados; proveen una forma fácil de manejar claves de longitud variable; algunos ahorran espacio guardando parte de la clave dentro de la estructura de búsqueda; y proveen un acceso muy rápido a los datos. Las desventajas son que los datos sesgados pueden llevar a un árbol degenerado con una mala eficiencia, especialmente porque los datos compuestos por caracteres son sesgados, pudiendo ser muy ineficientes en el uso del espacio.

Vamos a examinar varios métodos, donde cada uno corrige un problema inherente del anterior, llegando a un método importante que es bastante útil en aplicaciones de búsqueda donde tenemos claves muy largas. Uno de los problemas de hacer radix sorting es que no podemos lidiar fácilmente con claves iguales (se podría usar una lista de punteros a los distintos registros). En general, asumimos que todas las claves con las que vamos a trabajar son distintas.

3.4.1. Árboles de búsqueda digital

El método de radix searching más simple son los árboles de búsqueda digital. Los Árboles de Búsqueda Digital son parecidos a los ABB en cuanto a que para guardar una clave, vamos recorriendo el camino en el árbol que nos lleva a su posición. La diferencia está en que la posición no está determinada por

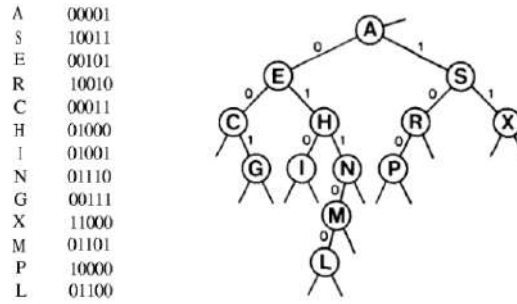


Figura 3.4: Árboles de búsqueda digital.

comparaciones de los valores completo de la clave, sino que por comparaciones sobre los bits (dígitos o caracteres) de la clave. En el primer nivel se utiliza el bit más significativo, en el segundo nivel el segundo más significativo, y así siguiendo hasta encontrar un nodo hoja.

En la Fig. 3.4 podemos ver un ejemplo de un Árbol de Búsqueda Digital. Supongamos que queremos insertar el elemento X en el árbol. Para ello, en cada paso nos movemos a derecha o a izquierda según si la i -ésima posición de la clave tenemos un 1 o un 0. Una vez alcanzamos una hoja, simplemente se coloca a X como hijo izquierdo o derecho de dicha hoja, según corresponda.

Podemos ver que en el peor caso, el costo de hacer búsquedas e inserciones va a depender del tamaño máximo de una clave, a diferencia de los ABBs que dependían de la cantidad de claves. Si lo miramos más en detalle, podemos notar que en realidad la altura del árbol de búsqueda digital depende del mayor número de bits sucesivos iguales entre dos claves cualesquiera, es decir, depende del mayor prefijo común entre dos claves cualesquiera. En el caso promedio, si los bits de los distintos elementos sean igualmente probables 1 o 0, podemos ver que a lo sumo se van a necesitar $O(\log n)$ comparaciones de claves completas. (Necesitamos hacer comparaciones de claves completas porque estamos almacenando el valor de la clave dentro de cada nodo, y a partir de su posición solo podemos deducir que la clave tiene cierto prefijo y no podemos deducir cuál es su valor.) Por lo tanto, resultan una alternativa atractiva frente a los ABB tradicionales, siempre que sea fácil realizar estas comparaciones de bits (este tipo de algoritmos son más fáciles de implementar en lenguajes con poca abstracción).

3.4.2. Tries

Las claves de búsqueda suelen ser muy largas, posiblemente de más de 20 caracteres. En estas situaciones, el costo de comparar una clave por igualdad puede terminar siendo dominante en el costo de la búsqueda, por lo que no puede ser ignorado. En los árboles digitales de búsqueda todavía teníamos que comparar claves en cada nodo. En esta sección vamos a ver una estructura de datos que nos permita implementar diccionarios comparando un solo dígito en cada nodo.

Para evitar tener que guardar en cada nodo una clave, lo que vamos a hacer es poner todas las claves en las hojas. Por lo tanto, vamos a tener dos tipos de nodos: los nodos internos, que simplemente contienen los enlaces a otros nodos, y los nodos hoja, que contienen las claves y no tienen enlaces. Esta estructura se conoce como *trie*, derivado de la palabra *retrieval*. Para buscar una clave en el árbol, vamos a movernos según los bits de la clave, pero no vamos a compararla con ninguna clave hasta llegar a una hoja. Cada clave en el árbol se almacena en una hoja que pertenece al camino descrito por el patrón de dígitos más significativos de la clave, y completamos la búsqueda comparando la clave deseada con la clave almacenada en la hoja. Es una estructura apropiada cuando tenemos muchas palabras que empiezan con la misma secuencia de caracteres, es decir, cuando la cantidad de prefijos distintos entre todas las palabras en el conjunto es mucho menor que la longitud total de todas las palabras.

En general, un trie es un árbol k -ario para alfabetos de $k - 1$ elementos posibles. En este tipo de árboles, los ejes representan componentes de las claves, por ejemplo, si las claves son strings, los ejes representan caracteres. Se tiene que cada camino desde la raíz hasta una hoja se corresponde con una palabra en el conjunto representado. De esta manera, los nodos del trie se corresponden con los prefijos de las palabras del conjunto representado.

En principio, un trie admite claves de cualquier longitud, por lo que es necesario tener un dígito adicional que nos permita saber que la clave terminó ahí ("§"). Si las claves fueran de una longitud prefijada, no sería necesario tener este dígito adicional, ya que la terminación de una clave vendría dada por el hecho de encontrarnos en el nivel k . Luego, el trie pasaría a ser un árbol k -ario para alfabetos de k elementos. También es posible tener distintos tipos de nodos con distinta aridad, por ejemplo, podríamos tener nodos con una alta aridad en los primeros niveles y nodos con baja aridad en los últimos (este método se dice *híbrido*).

Ejemplo

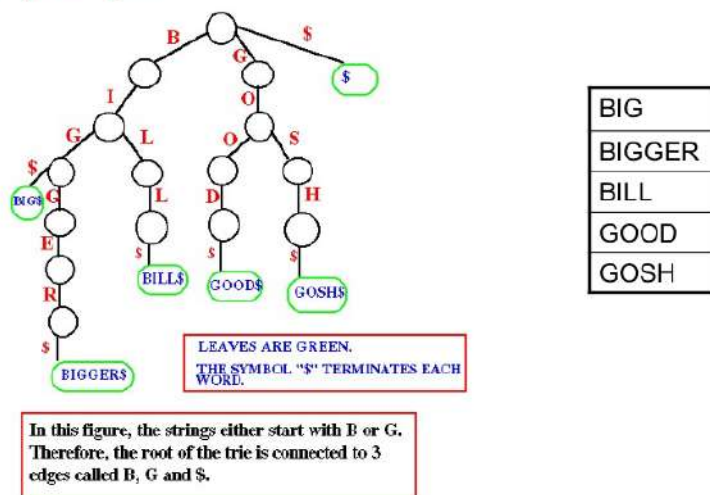


Figura 3.5: Ejemplo de Trie.

Para un trie construido a partir de n claves de b bits, bajo el modelo uniforme, se puede probar que se necesitan $\Theta(\log n)$ comparaciones **de bits** en el caso promedio, asumiendo una distribución uniforme de la codificación de los elementos², y en el peor caso $O(b)$ comparaciones de bits. También vale la pena mencionar que la estructura del trie es independiente del orden en el que se insertan las claves, es decir, hay un **único** trie para cada conjunto de claves (cosa que no pasaba ni en los ABB ni en los ABD).

El problema que tenemos en los Tries clásicos es la ramificación que surge al momento de insertar dos claves que tienen un prefijo muy largo en común. Por ejemplo, las claves que difieren en tan solo el último bit requieren de un camino cuya longitud sea igual a la longitud de la clave, sin importar la cantidad de claves en el árbol. Por lo tanto, la cantidad de nodos internos puede ser mayor que la cantidad de claves.

Implementación de Tries : Básicamente vamos a tener los dos formatos generales para representar varias variables: arreglos (memoria estática) y listas (memoria dinámica). La primera posibilidad sería tener en cada nodo interno un arreglo de punteros de longitud $k + 1$.

²dada una posición i , es igualmente probable encontrar un 1 que un 0

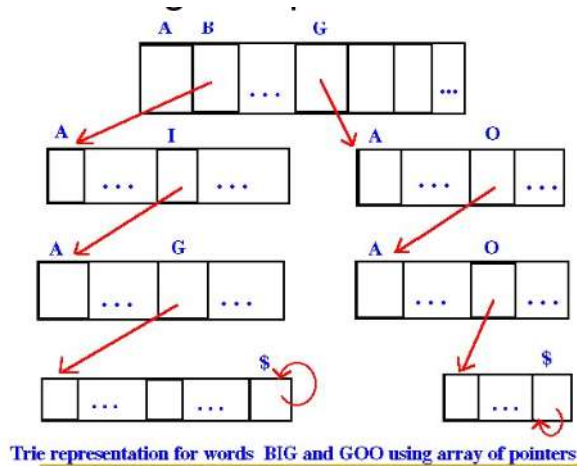


Figura 3.6: Implementación de Trie con Arreglos

En principio, en cada nodo estaríamos utilizando un arreglo de tamaño $k + 1$, siendo k la cantidad de elementos de nuestro alfabeto. El problema que tiene este esquema es posible que gran parte de las entradas no se utilicen, por lo que se estaría desperdiciando mucha memoria, especialmente cuando tenemos un alfabeto grande y un diccionario chico. La ventaja que tiene es que la eficiencia en términos de tiempo es extrema, al tener un acceso rápido a cada posición del arreglo.

La segunda posibilidad es tener una representación basada en listas, donde en cada nodo se almacenan todas las posibles continuaciones a través de una lista encadenada. Es decir, cada nodo va a tener un puntero a la lista de hijos, y además cada nodo va a tener un puntero hacia su hermano (formando una lista de hermanos).

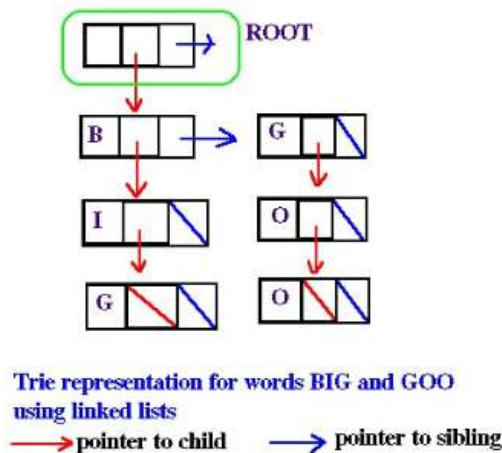


Figura 3.7: Implementación de Trie con Listas

Esta solución es eficiente en términos de tiempo solo si hay pocas claves (la lista de hermanos es corta), pero es mucho más eficiente en términos de memoria.

Optimizaciones : Para terminar, vamos a ver algunas variantes de los Tries que son utilizadas en la práctica. La primera optimización son los **Tries Compactos**, en los que colapsamos las cadenas que llevan a una única hoja. Es decir, si en algún momento llegamos a tener una serie de nodos en los que cada nodo tiene un único hijo, es innecesario representarlo con la cadena completa, ya que podemos directamente conectar al primer nodo de la cadena con el nodo hoja.

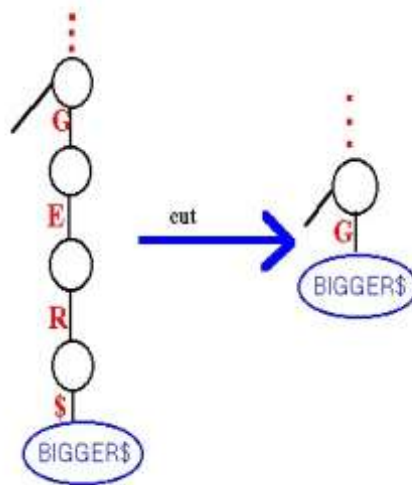


Figura 3.8: Tries Compactos

La segunda optimización que podemos hacer se basa en el hecho de que podemos compactar aquellos caminos intermedios que no tienen bifurcaciones. Esta segunda optimización nos lleva a la estructura **Patricia**, en donde un eje puede representar más de un dígito.

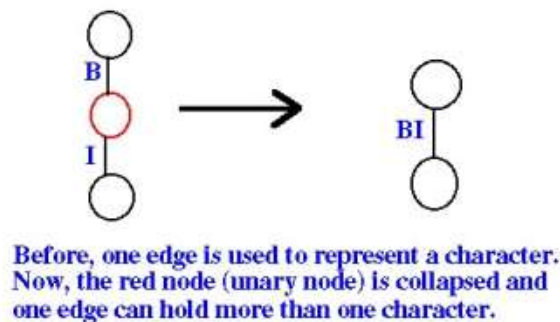


Figura 3.9: Tries Patricia

De esta manera, tenemos una mejora en el aprovechamiento del espacio, y se puede demostrar que la altura de un árbol Patricia binario está acotada por el número de claves. El problema que tiene esta estructura es que aumenta la dificultad de los algoritmos, al necesitar ahora de poder distinguir ejes que representen un solo dígito de aquellos que representen varios, trayendo un cierto overhead en los algoritmos. En conclusión, Patricia es el método de radix searching perfecto: consigue identificar los bits que distinguen las claves de búsqueda y los transforma en una estructura de datos (sin nodos sobrantes) que permite encontrar rápidamente cualquier clave en la estructura. Su atributo más útil es que puede usarse con claves de longitud variable de forma eficiente.

En todos los otros métodos de búsqueda que vimos, la longitud de las claves está incorporada de alguna manera en el algoritmo de búsqueda, por lo que el tiempo de ejecución es dependiente de la longitud de las claves al igual que el número de claves. Las ganancias que obtenemos va a depender del método de acceso a los distintos dígitos. Por ejemplo, si tenemos una computadora que puede acceder eficientemente a bytes de datos, y tenemos claves de 1000 bits, entonces Patricia tan solo requeriría acceder a 10 bytes para buscar la posible ubicación de la clave, más una única comparación de claves de 1000 bits. En cambio, los otros métodos requieren de varias comparaciones sobre las claves completas.

Formalmente, si tenemos claves de b bits (variable), bajo el modelo logarítmico, las operaciones del AVL dejan de ser $O(\log n)$ y pasan a ser $O(b \cdot \log n)$, ya que todas las operaciones involucradas son comparaciones simples sobre **claves** de b bits. En cambio, bajo el modelo logarítmico, las operaciones del trie siguen siendo $O(b)$ en el peor caso, ya que todas las operaciones involucradas son comparaciones sobre **bits** [2].

3.5. Hashing

Muchas aplicaciones requieren de un conjunto dinámico que dé soporte a las operaciones de diccionario de inserción, búsqueda y borrado. Por ejemplo, un compilador que traduce un lenguaje de programación mantiene una tabla de símbolos, en la cual las claves son cadenas de caracteres arbitrarias correspondientes a identificadores en el lenguaje. Una tabla hash es una estructura de datos efectiva para implementar diccionarios. A pesar de que la búsqueda de un elemento en una tabla hash puede tomar tiempo $\Theta(n)$ en el peor caso, en la práctica resultan extremadamente eficientes. Bajo asunciones razonables, el tiempo promedio de búsqueda de un elemento en una tabla hash es $O(1)$. Otro uso de las tablas hash es cuando trabajamos con el procesamiento de datos almacenados en memoria secundaria (ver *Hashing Extensible*).

Los arreglos trabajan con **direccionamiento directo**. Supongamos que tenemos que almacenar 10000 perfiles distintos, identificables a partir de una clave; en este caso, el número de DNI. El DNI es un número que va entre 0 y $10^8 - 1$, por lo que si quisiéramos utilizar un diccionario implementado a partir de un arreglo, tendríamos que tener 10^8 entradas del mismo. La ventaja de utilizar un arreglo es que podemos acceder a los datos del perfil en $O(1)$, pero estamos desperdiciando mucha memoria, ya que solo tenemos que almacenar 10000 perfiles distintos. En general, utilizar arreglos para representar conjuntos resulta eficiente cuando la cantidad de claves almacenadas es cercana a la cantidad total de claves posibles. Cuando esto no se cumple, las tablas hash resultan una alternativa eficiente, ya que estas utilizan un espacio proporcional a la cantidad de claves utilizadas. La idea es que en lugar de utilizar la clave como índice, el índice sea *computado* a partir de la clave. En esta sección vamos a ver cómo se calculan los índices, las distintas formas de direccionamiento y cómo se manejan las *colisiones* en las tablas hash.

Para indexar con otro tipos de datos, escapando de la necesidad de utilizar enteros, lo que vamos a hacer es encontrar una función de correspondencia que mapee a cada dato posible con un entero. La resolución de este problema se conoce como **Pre-hashing**. Una primera alternativa (teórica) consiste en aprovechar el hecho de que los datos en una computadora se representan como una tira de bits, la cual podemos interpretar como un entero. En la práctica, se cuenta con una función de pre-hasheo, que depende de la implementación de cada lenguaje de programación, la cual nos devuelve un entero a partir de la clave. Este tipo de funciones van a respetar (idealmente) que $\text{pre-hash}(x) = \text{pre-hash}(y) \iff x = y$. Una posible función de pre-hashing para string puede ser asociar a cada caracter su código ASCII, y al string asignarle el número obtenido en base 2, por ejemplo, "ppt" = $112 \cdot 2^2 + 112 \cdot 2 + 116 = 788$.

Cuando trabajamos con una tabla hash, vamos a representar un diccionario a partir de una tupla $\langle T, h \rangle$, donde T es un arreglo con m celdas, y $h : \text{Claves} \rightarrow \{0, \dots, n - 1\}$ es la función de hashing, que nos da una correspondencia entre el conjunto de claves posibles al conjunto de las posiciones de la tabla (estos índices son conocidos como *pseudoclaves*). De esta manera, la posición de un elemento en el arreglo se calcula a través de la función h .

Tabla hash

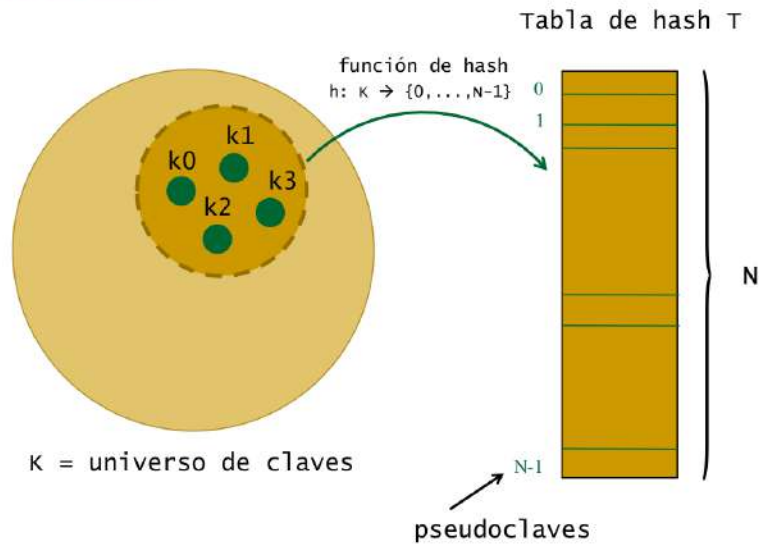


Figura 3.10: Tabla Hash

Existen muchos métodos para definir funciones hash:

- División: $h(k) = k \bmod m$. En general, se eligen un tamaño m para la tabla hash tal que sea un número primo no demasiado cercano a una potencia de 2.
- Partición: Se particiona la clave k en muchas claves k_1, k_2, \dots, k_n , y se toma como función de hash $h(k) = f(k_1, k_2, \dots, k_n)$.
- Extracción: Se usa solo una parte de la clave para calcular la posición, por ejemplo, si tenemos una clave de 12 cifras, se consideran solo las 6 centrales.
- Mid-square, etc.

Siempre se tiene como objetivo romper de alguna manera cualquier estructura o patrón previo que traigan consigo las claves, distribuyendo de forma uniforme las distintas claves.

Una función hash se dice **perfecta** si para todo par de claves distintas, las pseudoclaves asociadas son distintas. El problema que tenemos es que para poder asegurar que una función hash sea perfecta para cualquier conjunto de claves, se requiere que el tamaño de la tabla hash sea mayor o igual al cardinal del universo de claves, que es lo que queríamos evitar. Por lo tanto, en la práctica lo que se hace es tener una función hash que no sea perfecta, donde el tamaño de la tabla sea mucho menor a la cantidad de claves posibles.

En consecuencia, vamos a tener escenarios **frecuentes** en los que dos claves distintas van a ser mapeadas al mismo lugar, originando una **colisión**. Aún suponiendo una distribución uniforme entre las pseudoclaves, si elegimos 23 claves al azar dentro del universo de claves con una tabla hash de tamaño 365, la probabilidad de que dos de ellas colisionen es de más del 50% (Paradoja de los cumpleaños).

Como no podemos tener una función hash perfecta, lo que buscamos es una función que cumpla con la **uniformidad simple**, es decir, que las claves se distribuyan lo más equitativamente posible entre todas las posiciones del arreglo. Formalmente, una función hash cumple con la uniformidad simple si:

$$\sum_{k \in K: h(k)=j} P(k) \approx 1/m \quad \forall j,$$

siendo $P(k)$ la probabilidad de encontrarse con la clave k y m el tamaño de la tabla. El problema de construir este tipo de funciones es que, normalmente, no se conoce la probabilidad de aparición de las claves.

Entonces, se necesitan de métodos para la resolución de colisiones, los cuales se diferencian por la forma en la que ubican a los elementos involucrados en una colisión. Tenemos dos familias principales para la resolución de colisiones:

3.5.1. Concatenación

En el método de resolución de colisiones por concatenación, colocamos a todos los elementos que hashan a la misma posición de la tabla en una misma lista enlazada. Es decir, a la i -ésima posición de la tabla se coloca la cabeza de la lista que almacena a todos los elementos tales que $h(k) = i$; en caso de que no hayan elementos, se coloca un *NIL*. Es tentador intentar mejorar este esquema reemplazando estas listas de colisión por árboles balanceados, pero esto no vale la pena si el *factor de carga* se mantiene pequeño, a menos que sea esencial mejorar el costo asintótico en el peor caso.

En este caso, la inserción de un elemento k consiste en colocar al principio (o al final) de la lista asociada a la posición $h(k)$. En general, vamos a tener funciones hash que sean eficientes, es decir, que se resuelvan en tiempo constante para cualquier clave, por lo que la operación de inserción nos queda en $O(1)$. La operación de búsqueda consiste en recorrer linealmente la lista asociada a la posición $h(k)$, por lo que el costo de esta operación depende linealmente de la longitud de la lista asociada a $h(k)$. Para la operación de borrado, nuevamente recorreremos secuencialmente la lista asociada a $h(k)$, y borramos el elemento deseado, por lo que el costo nos queda lineal en función del tamaño de la lista.

Para analizar la longitud de las listas, definimos el **factor de carga** α para una tabla T con m posiciones que almacenan n elementos como $\alpha = n/m$, es decir, la longitud promedio de las listas. En el peor caso, todas las claves terminan hashadas en la misma posición, creando una lista de tamaño n . Por lo tanto, nos queda que el tiempo de cómputo de las operaciones en el peor caso es de $\Theta(n)$. La eficiencia del caso promedio va a depender de qué tan bien disperse las claves sobre la tabla la función de hash, en promedio. Suponiendo uniformidad simple en la función y que se puede computar la pseudoclave en tiempo constante, se puede probar que en el caso promedio tenemos un costo $\Theta(1 + \alpha)$.

En conclusión, si tenemos que la cantidad de posiciones en T es al menos proporcional a la cantidad de elementos, nos queda que $\alpha = n/m \approx 1$, por lo que las operaciones nos quedan en $O(1)$ en el caso promedio.

3.5.2. Direccionamiento abierto

En el método de resolución de colisiones por direccionamiento abierto, todos los elementos se guardan en la tabla. Es decir, cada entrada en la tabla contiene un elemento o *NIL*. Cuando hacemos una búsqueda, examinamos de forma sistemática las entradas de la tabla hasta que encontremos el elemento o que tengamos la certeza de que el elemento no está. No se almacenan listas ni elementos almacenados por fuera de la tabla, por lo que es posible que esta se llene ($\alpha \leq 1$). Evidentemente podríamos almacenar listas dentro de la tabla en aquellos lugares que actualmente estén disponibles, pero la ventaja del direccionamiento abierto es que no se requieren de punteros adicionales. En lugar de punteros, lo que *computamos* la secuencia de entradas a ser examinadas. La memoria adicional que ganamos por no usar punteros provee a la tabla un mayor número de entradas por la misma cantidad de memoria, potencialmente reduciendo la cantidad de colisiones y una búsqueda más rápida.

Para realizar una inserción usando direccionamiento abierto, vamos probando de forma sucesiva distintas casillas de la tabla hasta encontrar alguna que esté vacía. Por lo tanto, la función de hash no solo va a depender de la clave, sino que también de la cantidad de intentos que se hicieron para insertar este elemento, es decir, $h(k, i)$ nos define la posición para la clave k en el i -ésimo intento. Cuando ocurre una colisión, lo que hacemos es aumentar i , y probar devuelta con $h(k, i + 1)$ para que, eventualmente, encontrar un espacio vacío. Es necesario que para cada clave k , la secuencia de pruebas

$$\langle h(k, 0), h(k, 1), \dots, h(k, m - 1) \rangle$$

sea una permutación de $\langle 0, 1, \dots, m - 1 \rangle$, para que cada entrada en la tabla sea eventualmente considerada como casilla para cada clave nueva.

Decimos que una función de hash cumple con la propiedad de **uniformidad** si la secuencia de pruebas de cada clave es igualmente probable a cualquiera de las $m!$ permutaciones de $\langle 0, 1, \dots, m - 1 \rangle$.

Esta propiedad generaliza la noción de uniformidad simple que habíamos definido anteriormente para funciones de hash que generan secuencias de pruebas. Implementar funciones de hash que realmente satisfagan la propiedad de uniformidad es demasiado complicado, sin embargo, en la práctica se utilizan aproximaciones como el *hashing doble*.

El algoritmo de búsqueda consiste en ir visitando las distintas posiciones dadas por la función de hash, hasta encontrar alguna en la que esté el elemento o que se encuentre vacía. Notemos que estamos asumiendo que las claves nunca son borradas de la tabla. Cuando permitimos el borrado de una clave de la posición i , no podemos simplemente marcarla como *NIL* en su lugar. Si lo hiciéramos, es posible que seamos incapaces de recuperar cualquier clave que haya pasado por esa posición y la haya encontrado ocupada. Podemos resolver este problema marcando esa entrada como *borrada* en lugar de *NIL*. Luego, tendríamos que modificar el algoritmo de inserción para que trate ese tipo de entradas como si estuvieran vacías de manera tal que podamos insertar una nueva clave. No debemos modificar el algoritmo de búsqueda, ya que pasará por encima de las entradas marcadas como *borrada*. Sin embargo, cuando usamos el valor de *borrado*, el tiempo de búsqueda deja de depender en el factor de carga $\alpha \leq 1$, y es por este motivo que el método de concatenación suele ser más utilizado cuando se necesita la operación de borrado.

Vamos a diferenciar los métodos de direccionamiento abierto a partir de la **técnica de barrido** que utilicen, es decir, según cómo la función $h(k, i)$ va recorriendo todas las posiciones de la tabla ante una colisión. Las tres más típicas son el Barrido Lineal, el Barrido Cuadrático y el Hashing doble. Hacemos esta diferenciación porque cambia la complejidad de cálculo y su comportamiento respecto a los fenómenos de **aglomeración** (cómo se distribuyen las clave cuando ocurren colisiones).

Barrido Lineal : En este caso, la función de hash $h(k, i)$ va a ser igual a una función de hashing $(h'(k) + i) \bmod m$, para $i = 0, 1, \dots, m - 1$. Dada una clave k , primero probamos con $h'(k)$, la casilla dada por la función de hash auxiliar. Luego probamos con $h'(k) + 1$, y así siguiendo hasta $m - 1$. Después damos la vuelta a la tabla y seguimos con $0, 1, \dots, h'(k) - 1$. Como $h'(k)$ determina completamente la secuencia de pruebas, solo tenemos m secuencias distintas.

La técnica de barrido lineal es fácil de implementar, pero tiene el problema de **aglomeración primaria**. Se empiezan a armar segmentos largos de casillas todas ocupadas, incrementando el tiempo promedio de búsqueda. Estas aglomeraciones aparecen porque una casilla vacía precedida por i entradas llenas tiene probabilidad $(i + 1)/m$ de llenarse, porque la secuencia de pruebas de cualquier clave que caiga en cualquiera de las i entradas va a terminar en esta casilla. Esto lleva a que cada vez sea más probable caer en estos segmentos, extendiéndolos cada vez más.

Por ejemplo, en la Fig. 3.11 podemos ver que tenemos la secuencia de inserciones: 2, 103, 104, 105, ... Primero se inserta el 2, como es el primer intento $i = 0$, por lo que $h(2, 0) = (h'(2) + 0) \bmod 101$, siendo $h'(k) = k \bmod 101$, el 2 va a la posición 2. Luego, si inserta el 103, como es el primer intento $i = 0$, por lo que $h(103, 0) = 103 \bmod 101 = 2$, pero la posición 2 está ocupada. Entonces, se vuelve a intentar insertar con $i = 1$, y esta vez se coloca en la posición 3, que estaba libre.

- $h(k, i) = (h'(k)+i) \bmod 101$
- $h'(k)=k \bmod 101$
- Secuencia de inserciones {2, 103, 104, 105,...}
- Caso extremo, pero el problema existe!

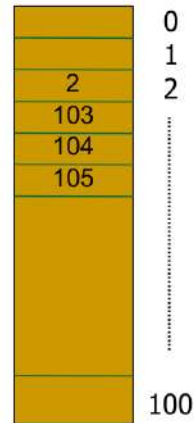


Figura 3.11: Ejemplo de Aglomeración Primaria.

Barrido Cuadrático : En este caso, la función de hash es de la forma

$$h(k, i) = (h'(k) + c_1 \cdot i + c_2 \cdot i^2) \bmod m,$$

donde c_1 y c_2 son constantes. La posición inicial es $h'(k)$, y las posiciones siguientes son desplazamientos que dependen de forma cuadrática de la cantidad de intentos. De esta manera, se evitan los grandes segmentos de la aglomeración primaria, pero para hacer uso de toda la tabla, los valores de c_1, c_2 y m están limitados. El problema que seguimos teniendo es que para cualquier par de claves k_1, k_2 tales que $h'(k_1) = h'(k_2)$, toda su secuencia de pruebas es exactamente la misma. Esta propiedad tiene un efecto más suave de aglomeración, conocida como **aglomeración secundaria**. Además, al igual que en la técnica de barrido lineal, $h'(k)$ determina la secuencia de pruebas completa, por lo que solo hay m secuencias distintas.

Hashing Doble : Ofrece uno de los mejores métodos disponibles para direccionamiento abierto. La idea es que el barrido también depende de la clave, utilizando una segunda función de hashing:

$$h(k, i) = (h_1(k) + i \cdot h_2(k)) \bmod m.$$

Por lo tanto, a diferencia de los casos de barrido lineal y de barrido cuadrático, la secuencia de pruebas depende de dos formas de la clave, ya que tanto la posición inicial como el desplazamiento varían. El valor de $h_2(k)$ debe ser coprimo con el tamaño de la tabla m para que la tabla pueda ser recorrida completamente. Una manera conveniente de asegurar esto es que m sea una potencia de 2 y que h_2 sea diseñada para que siempre devuelva un número impar. Cuando esto se cumple, la técnica de hashing doble termina generando $\Theta(m^2)$ secuencias de prueba, en lugar de las $\Theta(m)$ de las técnicas anteriores, ya que cada par $\langle h_1(k), h_2(k) \rangle$ define una secuencia distinta. En conclusión, la técnica de hashing doble no tiene aglomeración primaria, pero sí puede tener aglomeración secundaria, aunque esta última es menos probable que en los casos anteriores.

3.6. Colas de Prioridad

La Cola de Prioridad es un TAD parecido al TAD COLA, pero con la diferencia de que el orden en el que salen los elementos no está dado por el orden de llegada, sino que está dado por un atributo de los elementos, denominado **prioridad**, que normalmente está dada por un entero. Se permite el agregado de elementos en cualquier orden, y tenemos una operación de *desencolar* elemento, que nos devuelve el elemento encolado de mayor prioridad. Si tenemos empates, de alguna manera tenemos que resolverlo (por ejemplo, agregando un segundo orden de prioridad del tipo FIFO).

Para implementar este TAD, vamos a utilizar una estructura adecuada que tiene la misma complejidad en el peor caso que los árboles AVL. No vamos a obtener una mejora en la eficiencia en la implementación, pero vamos a tener una estructura mucho más simple que consigue implementar de

forma eficiente a las colas de prioridad, conocida como **Heap**. Además, al tratarse de estructuras mucho más sencillas, podemos pensar que el factor constante que multiplica a las operaciones de los heaps va a ser mucho más chico que el que acompaña a las de los AVL. Esto toma particular importancia cuando trabajamos con operaciones que toman tiempo logarítmico, porque si tenemos que dos algoritmos A, B cuestan $O(\log n)$, pero el factor constante de una es el doble de la otra $T_A(n) = 2T_B(n)$, entonces el tiempo que tarda A en procesar una entrada de tamaño n es similar al tiempo que tarda B en procesar una entrada de tamaño n^2 , dado que $2 \log n = \log n^2$.

Un *heap* es un árbol enraizado *perfectamente balanceado* que puede ser implementado eficientemente en un arreglo sin punteros explícitos. Cada nodo del árbol se corresponde con un elemento en el arreglo. El árbol está completamente lleno en todos los niveles exceptuando posiblemente el último, que está lleno desde la izquierda hasta algún punto. Esta estructura tiene numerosas aplicaciones, incluyendo la técnica de ordenamiento de *heapsort*. También puede ser usada para la representación eficiente de colas de prioridad dinámicas, como la lista de tareas a ser programadas por un sistema operativo. Es una estructura de datos ideal para encontrar el máximo elemento de un conjunto, removerlo, agregar un nuevo nodo, o modificarlo.

La idea es tener a los elementos parcialmente ordenados, permitiendo extraer el elemento más chico (o el más grande) y que al momento de sacarlo, es fácil reconstruir el invariante de la estructura, siendo el invariante de representación de un heap (también llamado **condición de heap**):

1. Es un Árbol Binario perfectamente balanceado (la longitud de todas las ramas difiere en a lo sumo 1).
2. La prioridad (posiblemente con una clave) de cada nodo es mayor que la de sus hijos.
3. Todo subárbol es un Heap.
4. El último nivel está lleno desde la izquierda (es *izquierdista*), o sea, todos los espacios libres en el último nivel tienen que estar a la derecha de los elementos de ese último nivel.

Cuando tenemos que a mayor clave, mayor prioridad, decimos que se trata de un **max-heap**. Cuando tenemos que a menor clave, mayor prioridad, decimos que tenemos un **min-heap**. De esta manera, cuando tenemos un *max-heap*, en la raíz vamos a encontrar el elemento más grande de toda la estructura, mientras que en un *min-heap*, vamos a encontrar el elemento más chico. Esto nos está diciendo que podemos **acceder** al próximo elemento en la cola de prioridad a ser desencolado, en tiempo constante $O(1)$, ya que lo podemos acceder directamente desde la raíz.

Las operaciones básicas definidas en el TAD COLA DE PRIORIDAD son:

- Vacía: Crea un heap vacío. En cualquiera de las implementaciones, tenemos costo $O(1)$.
- Próximo: Devuelve el elemento de máxima prioridad, sin modificar el heap. Consiste en mirar la raíz del heap, por lo que tenemos costo $O(1)$.
- Encolar: Agrega un nuevo elemento, teniendo que restablecer el invariante.
- Desencolar: Elimina el elemento de máxima prioridad, teniendo que restablecer el invariante.

Para implementar los heaps, nos sirven todas las representaciones usadas para árboles binarios: representación con punteros, donde contamos con nodos que tienen punteros explícitos a sus hijos; y la representación con arreglos, que es una representación implícita porque la posición en la que se encuentra cada elemento establece la relación que tiene con el resto de los nodos, siendo más eficiente que la representación explícita en cuanto al uso de memoria. En la representación con arreglos, vamos a tener dos longitudes: la longitud del arreglo, y el tamaño del heap (menor a la longitud del arreglo).

Para la representación con arreglos, lo que vamos a hacer es almacenar cada nodo v en la posición $p(v)$, donde $p(v)$ es una función recursiva que definimos de la siguiente manera:

- Si v es la raíz, entonces $p(v) = 0$.

- Si v es hijo izquierdo de u , entonces $p(v) = 2p(u) + 1$.
- Si v es hijo derecho de u , entonces $p(v) = 2p(u) + 2$.

En general, las ventajas de la representación con arreglos son la eficiencia en términos del espacio utilizado y la facilidad de navegación, o sea, al hecho de que a partir de la posición del padre, podemos determinar la posición de ambos hijos a partir de la función recursiva, y además podemos volver hacia atrás desde un hijo, a partir de la siguiente cuenta: posición del padre = $\lfloor (i - 1)/2 \rfloor$, siendo i la posición del hijo en el arreglo. Otras ventajas que tiene la implementación con arreglos es que tenemos una mejor vecindad espacial para un uso más eficiente de la memoria cache [1], y que permite una implementación sencilla de la operación DECREASEKEY en $O(\log n)$ (requiere de conocer de antemano el índice en el arreglo de la clave a decrementar).

La desventaja que tenemos es propia del uso de memoria estática, en la que nos vemos obligados muchas veces a redimensionar el arreglo a medida que se van agregando/eliminando elementos. Notemos que el motivo por el que podemos utilizar arreglos para representar heaps de forma eficiente sin desperdiciar memoria es que estamos trabajando con **árboles completos**, a diferencia de los árboles ABB y los árboles AVL en donde no podíamos garantizar esta propiedad (por lo que usar un arreglo para su representación podría significar un gasto excesivo de memoria).

Encolar : La operación de *encolar* consiste en dos pasos. Primero, insertamos el elemento a agregar en la primera posición disponible. Como sabemos que un heap es izquierdista, esta posición la podemos recordar y acceder vía un puntero al lugar que se ocupó último (no necesariamente coincide con el último elemento agregado). El siguiente paso consiste en recuperar el invariante del heap, es decir, la prioridad de cada nodo es mayor a la de sus hijos. En este caso, insertamos a un elemento en el último nivel sin tener en cuenta su prioridad, por lo que debemos *subirlo* hasta el nivel que corresponda. Esto último lo hacemos comparando únicamente al nodo con su padre. Mientras tenga mayor prioridad que su padre, lo intercambiamos de posición con su padre, siendo el peor caso que suba todos los niveles hasta volverse la raíz del heap. Por lo tanto, a lo sumo vamos a necesitar de $O(\log n)$ intercambios.

Desencolar : La operación de *desencolar* consiste en extraer el elemento de mayor prioridad del heap, es decir, eliminar la raíz. Nuevamente, consta de dos pasos. Primero, reemplazamos la raíz con la última hoja, y eliminamos la última hoja. Esto se debe a que un heap debe ser un árbol completo e izquierdista, por lo que si quitamos cualquier elemento (en este caso la raíz), una vez recuperado el invariante, no podemos tener esta última hoja. Luego, en la raíz nos quedó el elemento que se encontraba en la última hoja, por lo que no se estaría cumpliendo la propiedad de que para todo nodo, sus hijos tienen menor prioridad. Luego, lo que vamos a hacer es *bajar* este nodo que nos quedó en la raíz, hasta un lugar en el que se cumpla el invariante. Mientras este nodo p tenga al menos un hijo con mayor prioridad, lo que vamos a hacer es intercambiarlo con su hijo de mayor prioridad, siendo el peor caso que baje hasta volverse una hoja. Por lo tanto, a lo sumo vamos a necesitar de $O(\log n)$ intercambios.

Algoritmo de Floyd : Lo que queremos hacer es, dado un **arreglo** cualquiera, transformarlo en un heap (implementado como arreglo). La idea del algoritmo de Floyd es transformando localmente, de abajo hacia arriba (estrategia *bottom-up*), cada uno de los subárboles de la estructura hasta llegar a la raíz, de manera que al finalizar nos cumplimos con el invariante del heap. La idea es empezar con los subárboles más chicos que contengan hojas³, e ir aplicando la operación de *bajar* para convertirlos en heaps. Luego, una vez tenemos todos estos subárboles de altura k transformados en heaps, lo que hacemos es considerar los subárboles de altura $k + 1$, y se aprovecha el hecho de que todos los subárboles de altura k son heaps, por lo que a lo sumo se necesita de *bajar* un nodo por cada subárbol de altura $k + 1$. Se puede probar que el costo total nos queda en $O(n)$.

³Notemos que no hace falta hacer nada en el caso de considerar el subárbol que contenga solo una hoja, porque este ya cumple el invariante de heap.

Capítulo 4

Sorting

El problema de ordenamiento es uno de los problemas más clásicos y útiles de la informática. Vamos a dividir este tema en dos partes: ordenamiento interno y ordenamiento externo. El ordenamiento interno se realiza en la memoria principal de una computadora, donde podemos usar la capacidad de acceso arbitrario de la memoria de forma ventajosa. El ordenamiento externo es necesario cuando el número de objetos a ser ordenados es demasiado grande como para entrar en la memoria principal. En este contexto, el movimiento de datos entre la memoria principal y la memoria secundaria suele ser un cuello de botella, y debemos mover los datos en bloques para ser eficientes. El hecho de que es más conveniente mover de a bloques a los datos físicamente contiguos nos va a presentar una restricción respecto a los algoritmos de ordenamiento externo que podemos usar. El problema de ordenamiento externo será abarcado en una sección posterior.

Los algoritmos más simples de ordenamiento suelen tomar costo $O(n^2)$ para ordenar n objetos y son útiles solo para listas cortas (insertion sort, selection sort, bubblesort, shellsort, etc.). Uno de los algoritmos más populares de ordenamiento es *quicksort*, que toma tiempo $O(n \log n)$ en el caso promedio. Quicksort funciona bien en las aplicaciones comunes, sin embargo, en el peor caso tiene costo $O(n^2)$. Existen otros métodos, como *heapsort* y *mergesort*, que toman tiempo $O(n \log n)$ en el peor caso, pero su tiempo promedio puede que no sea tan bueno como quicksort. Por otro lado, mergesort es un algoritmo apropiado para ordenamiento externo (ver Sort-Merge). También vamos a considerar otros algoritmos llamados "bin" o "bucket" sort. Estos algoritmos solo funcionan en tipos especiales de datos, como los enteros que están limitados en cierto rango, pero cuando son aplicables, son extremadamente eficientes, tomando solo costo $O(n)$ en el peor caso.

Vamos a asumir que los objetos a ser ordenados son registros de uno o más campos. Uno de los campos, llamado la *clave*, es el tipo que tiene una relación de orden lineal \leq definida. Enteros, reales y arreglos de caracteres son ejemplos típicos de claves. El problema de ordenamiento consiste en colocar una secuencia de registros de manera tal que los valores de sus claves se encuentren en forma no decreciente. Existen varios criterios para evaluar el tiempo de ejecución de un algoritmo de ordenamiento interno. Los más comunes son medir la cantidad de pasos que se requieren para ordenar n registros, la cantidad de comparaciones entre pares de claves (importante si las claves son largas) y, si los registros son muy grandes, es de interés conocer la cantidad de veces en las que se debe mover a un registro.

Decimos que un algoritmo de ordenamiento es **estable** si dos registros i, j con claves iguales mantienen su orden relativo una vez ordenado el arreglo. Por ejemplo, vemos en la Fig. 4.1 tenemos un arreglo con cierto orden, en el que aparece primero i y luego j . Después de ordenar con un algoritmo *estable*, lo que se asegura es que i va a terminar antes que j en el arreglo. En cambio, después de ordenar con un algoritmo *inestable*, no tenemos esta garantía, sino que podría pasar que j termine antes que i , o también podría pasar que i termine antes que j en el arreglo. Esto resulta de interés cuando tenemos claves múltiples, por ejemplo una tupla \langle clave primaria, clave secundaria \rangle , en donde posiblemente nos gustaría poder ordenar por la clave primaria, y para desempatar utilizar la clave secundaria.

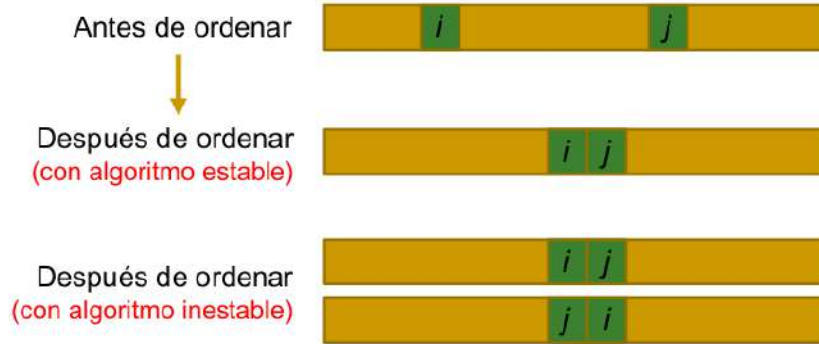


Figura 4.1: Estabilidad

4.1. Heap Sort

La idea es básicamente aplicar Selection Sort, es decir, seleccionar el mínimo elemento del sub-arreglo que ocupa las posiciones $i \cdot n - 1$, desde $i = 0, \dots, n - 2$, siendo n la cantidad de elementos del arreglo. La diferencia va a estar en que vamos a transformar el arreglo en un max-heap, utilizando el algoritmo de Floyd en $O(n)$, y luego vamos a descolocar n veces el máximo. Para hacer el algoritmo *in-place*, es decir, poder reutilizar el mismo arreglo en el que vienen los datos, vamos a aprovechar el hecho de que al momento de ir descolocando el heap, nos van a quedar libres los últimos $n - i$ espacios, en donde podemos ir colocando a los elementos a medida que los vamos descolocando. El costo de este algoritmo es $O(n \cdot \log n)$. Este algoritmo **no** es estable.

Supongamos que tenemos el siguiente arreglo: $\{(10, a), (5, a), (5, b), (5, c), (4, a), (3, a), (2, b)\}$, y veamos cómo el algoritmo termina cambiando el orden relativo entre $(5, a)$, $(5, b)$ y $(5, c)$. Al aplicar el algoritmo de Floyd, como este arreglo ya es un heap, no hace ningún intercambio. Ahora, al momento de descolocar, primero se descoloca el $(10, a)$, y suponiendo que en caso de empate baja por la rama izquierda, nos queda como raíz a $(5, a)$. Ahora, $(5, a)$ tiene como hijo izquierdo a $(5, c)$ y como hijo derecho a $(5, b)$, por lo que si descolocamos, nos queda como raíz a $(5, c)$. Si seguimos descolocando, lo que nos termina quedando es el arreglo ordenado $\langle (10, a), (5, a), (5, c), (5, b), (4, a), (3, a), (2, a) \rangle$. El caso en el que se baja por la rama derecha en caso de empate es similar.

4.2. Merge Sort

Es un algoritmo basado en la técnica de *Divide & Conquer*, técnica que se basa en *dividir* un problema en problemas similares, pero de tamaño más chico, llamados **sub-problemas** (todos los sub-problemas deben ser del mismo tamaño), para luego *combinar* las soluciones de los sub-problemas y así obtener la solución al problema original.

En este caso, lo que vamos a hacer es dividir al arreglo a la mitad, y a cada mitad del arreglo la ordenamos de forma recursiva. En caso de llegar a un arreglo conformado por un solo elemento, este arreglo ya está ordenado, por lo que podemos terminar la recursión. Una vez tenemos ambas partes ordenadas, las unimos haciendo un *Merge* en $O(n)$.

Entrada: Arreglo desordenado de n elementos

```

1 if  $n < 2$  then
2   | El arreglo está ordenado.
3 else
4   | Dividir el arreglo en dos partes iguales de  $n/2$  elementos.
5   | Ordenar recursivamente ambas mitades.
6   | Combinar ambas mitades ya ordenadas en un único arreglo

```

En general, el costo del peor caso de este algoritmo nos queda:

$$T(n) = \begin{cases} O(1) & \text{si } n = 1 \\ 2T(n/2) + O(n - 1) & \text{caso contrario} \end{cases}$$

y considerando que a lo sumo vamos a tener $O(\log n)$ recursiones, por el hecho de que en cada paso dividimos a la mitad el tamaño del arreglo, se puede demostrar que el costo total nos queda en $\mathcal{O}(n \cdot \log n)$ ¹.

4.3. Quick Sort

La esencia de quicksort es ordenar un arreglo $A[1], \dots, A[n]$ a través de elegir algún valor v en el arreglo como elemento *pivote*, alrededor del cual se ordena el arreglo. Idealmente, el pivote estaría lo más cercano a la clave mediana del arreglo, es decir, que sea precedida por la aproximadamente la mitad de las claves y que le siga la otra mitad. Luego, permutamos los elementos del arreglo de manera tal que para algún j , todos los registros con claves menores que v aparezcan en $A[1], \dots, A[j]$, y que todos aquellos cuyas claves sean mayores o iguales a v aparezcan en $A[j + 1], \dots, A[n]$. Se continúa aplicando recursivamente quicksort sobre ambos sub-arreglos. Como todas las claves del primer grupo son menores que todas las claves del segundo grupo, el arreglo termina ordenado. Notemos que la recursión termina cuando tenemos un arreglo un solo elemento (o cuando todas las claves son iguales), y que el algoritmo es *in-place*.

El problema que surge es cómo hacemos para encontrar al elemento mediano. Si bien existen algoritmos lineales que permiten encontrar el elemento mediano en un arreglo, estos vienen acompañados por una constante lo suficientemente alta como para que no nos resulte conveniente utilizarlos en la práctica. Luego, la idea va a ser tomar un elemento al que llamaremos *pivote*, y asumir que es lo suficientemente cercano al mediano.

Entrada: arreglo de n elementos

```
1 if array.length > 1 then
2   Elegir pivote
3   while Haya elementos en el array do
4     if elemento generico < pivote then
5       | insertar elemento en subarray1
6     else
7       | insertar elemento en subarray2
8   quicksort(subarray1)
9   quicksort(subarray2)
```

El costo de este algoritmo depende de la cantidad de comparaciones que se realicen a lo largo de las recursiones. En el peor caso, el pivote elegido siempre va a dividir el arreglo en un sub-arreglo de un solo elemento y en otro con los $n - 1$ elementos restantes, por lo que nos quedaría un costo de $O(n^2)$. En el mejor caso y en el caso promedio (asumiendo una distribución uniforme), se puede demostrar que nos queda un costo $O(n \cdot \log n)$.

En la práctica, la elección del pivote es fundamental, y se suelen tomar medidas para evitar caer en el peor caso. Algunos ejemplos pueden ser:

- permutar el input,
- elegir el pivote al azar,
- tomar k elementos aleatorios, obtener la mediana con algún algoritmo eficiente, y utilizarla como pivote,

¹Para esta demostración, se asume que n es potencia de 2. Esto se puede hacer ya que la función es **suave**, es decir, vale que $b \cdot f \in O(f)$ para todo $b \geq 2$. Si este no fuera el caso, no sería una cota exacta.

- llamar a los algoritmos simples de ordenamiento cuando los arreglos son lo suficientemente chicos como para que resulten más eficientes que quicksort (Knuth propone arreglos de tamaño 9),
- si las claves son grandes y costosas de mover, podemos usar punteros para no tener que mover las claves de lugar, a costa de un mayor uso de memoria y que acceder a las claves para realizar las comparaciones es más lento.

4.4. Bin Sorting

Más adelante vamos a mostrar que se necesitan de $\Omega(n \log n)$ pasos para ordenar n elementos, si no asumimos nada acerca de las claves. Sin embargo, si tenemos conocimiento de alguna propiedad de las claves, es posible que podamos ordenarlas en menor tiempo que $O(n \log n)$. Supongamos que queremos ordenar un arreglo de enteros A que sabemos que están en un rango $1, \dots, n$, sin duplicados, donde n es la cantidad de elementos. Entonces, podemos ordenar los elementos colocándolos en un arreglo auxiliar B según las claves:

Bin Sorting

```

1 for  $i = 1, \dots, n$  do
2    $B[A[i]] = A[i]$ 

```

Este código calcula a dónde pertenece el registro $A[i]$ y lo coloca en su lugar, en tiempo $O(n)$. Funciona correctamente cuando tenemos exactamente un registro con clave v para cada valor de v en $1, \dots, n$. Este es un ejemplo de "bin-sort", un proceso de ordenamiento en el que creamos un "bin" para que almacene todos los registros con cierto valor. En el caso general, vamos a tener que prepararnos para almacenar más de un registro en un bin y poder concatenarlos en el orden correcto. Para ello, vamos a tener un arreglo B de listas, que vendrían siendo los bins para cada tipo de clave. En general, si tenemos n elementos a ordenar y tenemos m claves distintas, el costo total de bin sorting nos queda en $O(n + m)$. Si tenemos que el número de claves $m \leq n$, entonces $O(n + m) = O(n)$.

Si $m = n^2$, el algoritmo anterior nos quedaría en $O(n^2)$. La pregunta que nos surge es, ¿se puede mejorar? La respuesta es que incluso si los posibles valores de las claves son $1, 2, \dots, n^k$, para cualquier k fijo, entonces existe una técnica que generaliza bin sorting que toma tiempo $O(n)$. Consideremos el caso en el que tenemos que ordenar n enteros en el rango de valores $0, \dots, n^2 - 1$. Vamos a ordenar a los n elementos en dos pasos. Primero usamos n bins, uno por cada entero $0, \dots, n - 1$. Colocamos cada entero i en el bin número $i \bmod n$, al final de la lista. Por ejemplo, si tenemos $n = 10$ y el arreglo $36, 9, 0, 25, 1, 49, 64, 16, 81, 4$, al finalizar el primer paso nos quedan 10 bins ordenados internamente en el mismo orden de aparición de la lista original. Luego, concatenamos los bins y nos queda $0, 1, 81, 64, 4, 25, 36, 16, 8, 49$.

Como segundo paso, lo que hacemos es volver a distribuir los números de la nueva lista, pero usando una estrategia de selección distinta. Ahora vamos a colocar el i -ésimo elemento en el bin $\lfloor i/n \rfloor$, respetando el orden de los elementos en la lista. Cuando concatenamos esta segunda lista, terminamos con una lista ordenada.

Para ver por qué este algoritmo funciona, nos tenemos que dar cuenta que cuando colocamos varios enteros en un bin, estos deben estar en orden creciente, ya que la lista resultante del primer paso los había ordenado por el dígito de más a la derecha (el menos significativo). Por lo tanto, en cualquier bin tenemos que los dígitos de más a la derecha forman una secuencia ordenada. De forma más general, podemos pensar a los enteros entre 0 y $n^2 - 1$ como números de dos dígitos en base n y usar el mismo argumento para ver que la estrategia funciona. Este algoritmo se conoce como *Radix Sorting*. La idea central detrás de radix sorting es hacer bin sort sobre todos los registros, primero en f_k , el dígito menos significativo, luego concatenar los bins, manteniendo los valores más chicos primero, volver a aplicar bin sort en f_{k-1} , y así hasta terminar con la lista ordenada. El costo total de radix sorting, para claves en el rango $0, \dots, n^k - 1$ es de $O(n + kn) = O(n)$. Si k aumenta a medida que aumenta n , por ejemplo si las claves son strings de longitud $\log n$, entonces nos quedaría $O(n \log n)$.

4.5. Árboles de decisión

Vamos a centrarnos en los algoritmos de ordenamiento que solo usan a los elementos para comparar dos claves. Podemos dibujar un árbol binario en el que cada nodo representa el *estado* del programa luego de realizar una cierta cantidad de comparaciones de claves. También podemos ver a un nodo como representante del orden inicial de los datos que traen al programa a este estado. Luego, un estado de un programa es esencialmente el conocimiento acerca del orden inicial del arreglo obtenido hasta el momento por el programa (la idea es que para poder ordenar un arreglo, necesitamos poder saber cuál es el orden inicial del mismo).

Si cualquier nodo representa dos o más posibles órdenes iniciales, entonces el programa no puede saber el ordenamiento correcto, por lo que debe realizar otra comparación de claves, como ¿es $k_1 < k_2$? Luego, podemos crear dos hijos para el nodo, donde el hijo izquierdo represente aquellos órdenes iniciales consistentes con que $k_1 < k_2$; mientras que el hijo derecho representaría aquellos órdenes que sean consistentes con el hecho de que $k_1 > k_2$. Entonces, cada hijo representa el estado consistente en la información acerca del padre más el hecho de que $k_1 < k_2$ o que $k_1 > k_2$, dependiendo si es el hijo izquierdo o el derecho.

En general, si tenemos que ordenar una lista de n elementos, existen $n!$ posibles salidas, que se corresponden con los posibles órdenes iniciales de la lista. Es decir, cualquiera de los n elementos podría ser el primero, cualquiera de los $n - 1$ restantes podría ser el segundo, etc. Por lo tanto, cualquier árbol de decisión que describa un correcto algoritmo de ordenamiento debe tener al menos $n!$ hojas. La longitud de un camino a una hoja nos da la cantidad de comparaciones que se necesitan para ordenar una lista. Luego, la longitud del camino más largo desde la raíz hasta una hoja es una cota inferior a la cantidad de pasos realizados por cualquier algoritmo de ordenamiento en el peor caso. Si un árbol tiene $n!$ hojas, sabemos que al menos tiene altura $\log n!$, por lo que cualquier algoritmo de ordenamiento debe ser $\Omega(\log(n!)) = \Omega(n \log n)$ (por la aproximación de Stirling) en el peor caso. Esto nos dice que tanto Merge Sort como Heap Sort tienen complejidad asintótica óptima (también vale para las operaciones de los árboles AVL). En el caso promedio, esta propiedad sigue valiendo, por lo que dichos algoritmos también son óptimos en el caso promedio.

Capítulo 5

Dividir y Conquistar

Ya vimos cómo *merge sort* nos servía como un ejemplo del paradigma de divide-and-conquer. Recordemos que la técnica algorítmica de DC consiste en resolver un problema de forma recursiva, aplicando tres pasos en cada nivel de la recursión. Primero, **dividir** el problema en un número de subproblemas que son instancias más chicas del mismo problema. Luego, **conquistar** los subproblemas resolviéndolos de forma recursiva. Si los subproblemas son lo suficientemente chicos, resolverlos de forma *ad hoc*. Para finalmente **combinar** las soluciones de los subproblemas en la solución del problema original (si existe algún costo asociado a la división del problema en subproblemas, se lo contamos en la etapa de combinación). Notemos que este método tiene sentido siempre y cuando la división y la combinación no sean operaciones excesivamente costosas.

Algoritmo DC

Entrada: Instancia genérica X

```
1 if  $X$  es lo suficientemente simple como para resolverlo ad hoc then
2    $\lfloor$  ADHOC( $X$ )
3 else
4   Descomponer  $X$  en sub-instancias  $X_1, X_2, \dots, X_k$ 
5   for  $i = 1 \dots k$  do
6      $\lfloor Y_i \leftarrow DC(X_i)$ 
7   Combinar las soluciones de  $Y_i$  para obtener una solución  $Y$  para la instancia original  $X$ .
```

Las recurrencias van de la mano con el paradigma de DC, porque nos dan una forma natural de caracterizar los tiempos de ejecución de este tipo de algoritmos. Una *recurrencia* es una ecuación o inecuación que describe la función en términos de su valor en entradas más chicas. Por ejemplo, vimos que para podíamos describir el costo de *merge sort* por la recurrencia:

$$T(n) = \begin{cases} \Theta(1) & \text{si } n = 1 \\ 2T(n/2) + \Theta(n) & \text{si } n > 1 \end{cases}$$

cuya solución dijimos (sin demostrar) que intuitivamente $T(n) \in \Theta(n \cdot \log n)$.

Las recurrencias pueden tomar muchas formas. Por ejemplo, un algoritmo recursivo podría dividir un problema en subproblemas de distintos tamaños (uno de tamaño $2/3$ y otro de tamaño $1/3$). El tamaño de los subproblemas no necesariamente está restringido a ser una fracción constante del tamaño del problema original. Por ejemplo, una versión recursiva de la búsqueda secuencial crearía un subproblema conteniendo solo un elemento menos que el problema original. Vamos a ver tres métodos para resolver recurrencias:

- El *método de sustitución*, en el que probamos una cota y usamos inducción para probar que es correcta.

- El *método del árbol de recurrencia*, que convierte la recurrencia en un árbol cuyos nodos representan los costos asociados a los niveles de la recursión. Usamos técnicas para acotar sumas para resolver la recurrencia.
- El *método maestro* que provee cotas para recurrencias de la forma

$$T(n) = \begin{cases} aT(n/b) + f(n) & \text{si } n > 1 \\ 1 & \text{si } n = 1 \end{cases}$$

con $a \geq 1$, $b > 1$ y $f(n)$ una función dada. Este último tipo de recurrencias aparecen frecuentemente, ya que caracterizan a los algoritmos DC que crean a subproblemas, cada uno de tamaño $1/b$ del tamaño del problema original, y que se pueden dividir y combinar en tiempo $f(n)$. Para ver otros métodos de resolución de ecuaciones de recurrencia, ver 4.7 de [5] y el método Akra-Bazzi.

5.1. Método de sustitución

El método de sustitución para resolver recurrencias se compone de dos pasos:

1. Suponer que la función tiene cierta forma.
2. Usar inducción para encontrar las constantes y mostrar que la solución es correcta.

Como ejemplo, determinemos una cota superior para la recurrencia

$$T(n) = 2T(n/2) + n.$$

Suponemos que la solución es $T(n) = O(n \log n)$. Debemos probar que $T(n) \leq c \cdot n \log n$ para una cierta constante $c > 0$. Empezamos asumiendo que vale para todos los positivos $m < n$, por lo que $T(n/2) \leq c \cdot n/2 \log(n/2)$. Si sustituimos en la recurrencia, nos queda que

$$\begin{aligned} T(n) &\leq 2T(n/2) + n \\ &= 2(c \cdot n/2 \log(n/2)) + n \\ &= cn \log(n/2) + n \\ &= cn \log(n) - cn \log(2) + n \\ &\leq cn \log(n) \quad \text{si } c \geq 1. \end{aligned}$$

Para completar la demostración, debemos tomar ventaja de que la notación asintótica nos pide probar que $T(n) \leq c \cdot n \log n$ para $n \leq n_0$, donde n_0 es una constante que podemos elegir. Luego, tomamos el caso $n_0 = 2$ y como toda recursión de $T(n)$ con $n \geq 2$ debe derivar en $T(2)$ o $T(3)$, basta con probar que existe un c tal que $T(2) \leq c \cdot 2 \log 2$ y $T(3) \leq c \cdot 3 \log 3$. Como $c = 2$ cumple, queda demostrado que $T(n) \in O(n \log n)$.

A veces podemos acertar en la cota asintótica de la solución de recurrencia, pero fallamos en demostrarla por inducción. El problema suele ser que estamos tomando una hipótesis inductiva demasiado débil como para probar la cota. Consideremos la recursión

$$T(n) = T(n/2) + T(n/2) + 1$$

Claramente, $T(n) \in O(n)$, por lo que intentamos probar que $T(n) \leq cn$. Si sustituimos en la recurrencia nos queda:

$$T(n) = cn/2 + cn/2 + 1 = cn + 1$$

que no implica que $T(n) \leq cn$ para cualquier elección de c . Lo que podemos hacer es *sustraer* un término de menor grado de nuestra solución. Es decir, vamos a probar que $T(n) \leq cn - d$, con $d \leq 0$ constante. Ahora nos queda que

$$T(n) \leq cn - 2d + 1$$

y tomando $d \leq 1$ quedaría demostrado que $T(n) \leq cn - d$. Notemos que no podíamos decir que $cn + 1 \in O(n) \Rightarrow T(n) \in O(n)$ porque estaríamos usando como argumento lo que queríamos demostrar: habíamos supuesto que $T(n) \leq cn$.

A veces, nos puede servir hacer un cambio de variables para facilitarnos la resolución de la recurrencia. Por ejemplo, si tenemos

$$T(n) = 2T(\sqrt{n}) + \log n$$

no se nos ocurre qué forma tiene a simple vista. Sin embargo, podemos simplificar esta recurrencia haciendo un cambio de variables, renombrando $m = \log n$. Entonces

$$T(2^m) = 2T(2^{m/2}) + m,$$

y si llamamos $S(m) = T(2^m)$ nos queda la siguiente recurrencia

$$S(m) = 2S(m/2) + m,$$

que ya sabemos resolver (es la de *merge sort*). Volviendo de $S(m)$ a $T(n)$, nos queda

$$T(n) = S(m) = O(m \log m) = O(\log n \cdot \log(\log n)).$$

5.2. Árbol de Recurrencia

Uno de los problemas que tiene el método de sustitución es que no conocemos ningún método general para conocer la forma que tiene la función. Afortunadamente, podemos usar los árboles de recurrencia para generar buenas aproximaciones. En un *árbol de recurrencia*, cada nodo representa el costo de algún subproblema en el conjunto de llamados recursivos. Podemos sumar los costos de cada nivel para obtener un conjunto de costos por nivel, y luego sumar a estos costos por nivel para obtener el costo total de la recursión. Si tenemos el suficiente cuidado con estas cuentas, podemos usarlo como solución de una función de recurrencia; sino, debemos formalizar mediante el método de sustitución. Veamos un ejemplo.

Queremos conocer una cota superior para $T(n) = 3T(n/4) + cn^2$. En la Fig. 5.1 podemos ver cómo derivar el árbol de recurrencia. La parte (a) de la figura muestra $T(n)$, que se expande en la parte (b) en un árbol equivalente representando la recurrencia. El término cn^2 en la raíz representa los costos asociados a la cima de la recursión (la combinación de los resultados de los subproblemas), y los tres subárboles de la raíz representan el costo asociado a resolver los subproblemas de tamaño $n/4$. En la parte (c) vemos el siguiente paso de este proceso, expandiendo cada nodo con costo $T(n/4)$ de la parte (b). El costo para cada uno de los tres hijos de la raíz es $c(n/4)^2$. Continuamos expandiendo cada nodo en el árbol siguiendo la recurrencia.

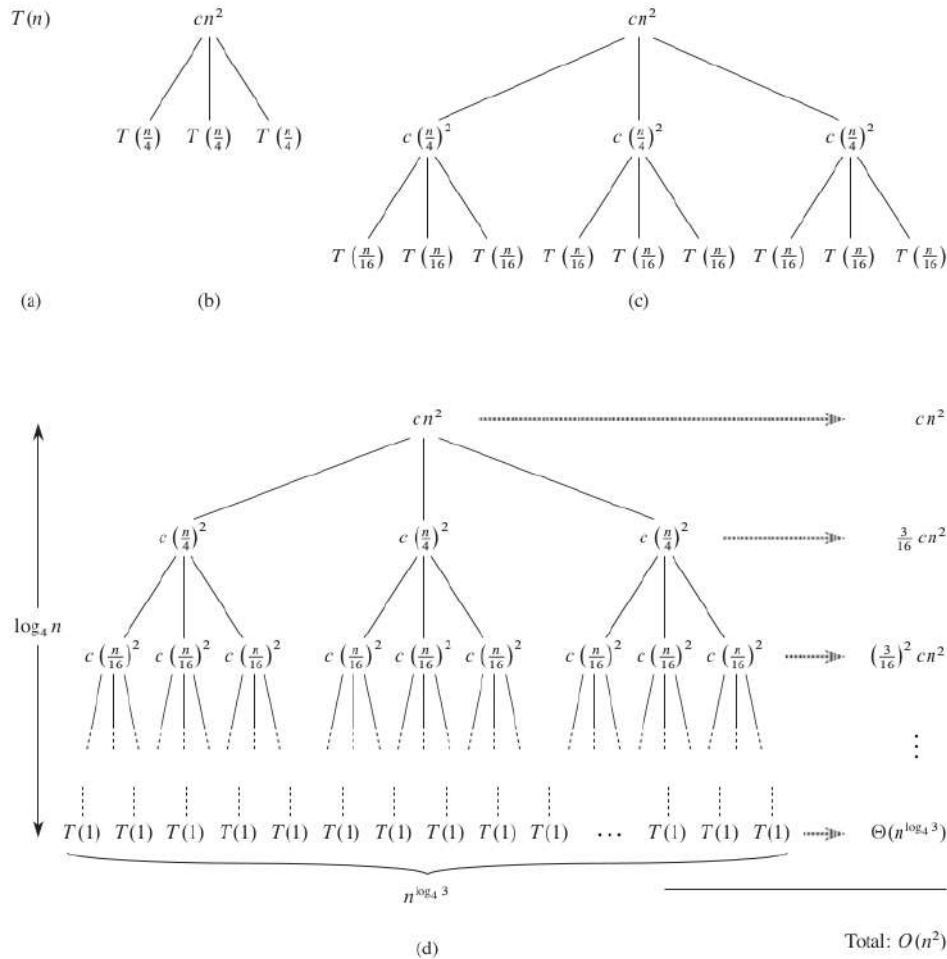


Figura 5.1: Construyendo un árbol de recurrencia para $T(n) = 3T(n/4) + cn^2$. Parte (a) muestra $T(n)$, que se va expandiendo progresivamente a (b)-(d) para formar el árbol de recurrencia. El árbol completamente expandido en la parte (d) tiene altura $\log_4 n$.

Como el tamaño de cada subproblema decrece por un factor de 4 en cada nivel, eventualmente vamos a llegar a un caso base que sabemos resolver *ad hoc*. Para simplificar, suponemos que el caso base es cuando $n = 1$, por lo que llegamos a este caso luego de $\log_4 n$ niveles. Por lo tanto, nos queda un árbol de altura $\log_4 n + 1$.

Ahora debemos determinar el costo de cada nivel del árbol. Sabemos que cada nivel tiene 3 veces más nodos que el nivel de arriba, por lo que la cantidad de nodos a profundidad i es 3^i . Por otro lado, como el tamaño de los subproblemas se reduce por un factor de 4 en cada nivel, tenemos que cada nodo a profundidad i tiene un costo asociado de $c(n/4^i)^2$. Multiplicando ambos resultados, podemos ver que el costo total de todos los nodos a profundidad i es $(3/16)^i cn^2$, exceptuando a las hojas. El último nivel, como está a profundidad $\log_4 n$, tiene $3^{\log_4 n} = n^{\log_4 3}$ nodos, y como cada nodo cuesta $T(1)$, nos queda que el último nivel cuesta $\Theta(n^{\log_4 3})$. Si sumamos todos los costos por nivel, nos queda que $T(n) = O(n^2)$.

5.3. Método Maestro

El método maestro nos da una receta para resolver recurrencias del tipo

$$T(n) = aT(n/b) + f(n),$$

donde $a \geq 1$ y $b > 1$ son constantes positivas y $f(n)$ es una función asintóticamente positiva. Para utilizar este método, vamos a tener que memorizar tres casos. Este tipo de recurrencias describen algoritmos que

dividen a un problema de tamaño n en a subproblemas de tamaño n/b que se resuelven en tiempo $T(n/b)$. La función $f(n)$ engloba todos los costos asociados a la división y combinación de los resultados de los subproblemas.

El método maestro depende del siguiente teorema.

Teorema Maestro: Sea $a \geq 1$ y $b > 1$ constantes, $f(n)$ una función, y sea $T(n)$ definida en los enteros no negativos por la recurrencia

$$T(n) = aT(n/b) + f(n).$$

Entonces, $T(n)$ tiene las siguientes cotas asintóticas:

1. Si $f(n) \in O(n^{\log_b(a)-\epsilon})$ para alguna constante $\epsilon > 0$, entonces $T(n) \in \Theta(n^{\log_b(a)})$.
2. Si $f(n) \in \Theta(n^{\log_b(a)})$, entonces $T(n) \in \Theta(n^{\log_b(a)} \cdot \log n)$.
3. Si $f(n) \in \Omega(n^{\log_b(a)+\epsilon})$ para alguna constante $\epsilon > 0$, y si para todo n suficientemente grande vale que $af(n/b) \leq cf(n)$ para alguna constante $c < 1$, entonces $T(n) \in \Theta(f(n))$.

En los tres casos, estamos comparando la función $f(n)$ con $n^{\log_b(a)}$. Intuitivamente, el mayor de las dos funciones determina la solución a la recurrencia. Por ejemplo, en el caso 1 la función $n^{\log_b(a)}$ es la mayor, y esto determina que $T(n) \in \Theta(n^{\log_b(a)})$.

Notemos que en el primer caso, $f(n)$ no solo debe ser menor a $n^{\log_b(a)}$, sino que debe ser *polinomialmente* menor. Esto quiere decir que $f(n)$ debe ser asintóticamente menor que $n^{\log_b(a)}$ por un factor n^ϵ para alguna constante $\epsilon > 0$. Esto quiere decir que no se están cubriendo todas las posibilidades para $f(n)$. Existe una brecha entre los casos 1 y 2 cuando $f(n)$ es menor que $n^{\log_b(a)}$, pero no polinomialmente menor. En el tercer caso, además, no solo $f(n)$ debe ser polinomialmente mayor a $n^{\log_b(a)}$, sino que debe satisfacerse la condición de que $af(n/b) \leq cf(n)$. Esta condición suele satisfacerse en las funciones polinomialmente acotadas. Si alguna de estas condiciones no se cumplen, no podemos usar el método maestro para resolver la recurrencia.

Capítulo 6

Memoria Secundaria

En este capítulo vamos a ver cómo se resuelven los problemas de *Sorting* y de *Searching* en **memoria secundaria**. Hasta ahora, solo nos interesaba hablar de complejidad asintótica y de TADs. Hay situaciones en las que la abstracción que hacemos para definir los problemas necesita un enfoque práctico que tenga en cuenta las propiedades físicas de una computadora.

Muchas aplicaciones de ordenamiento involucran el procesamiento de archivos muy grandes, demasiado grandes como para poder entrar en la memoria principal de cualquier computadora. Los métodos apropiados para este tipo de aplicaciones son llamados métodos *externos*, ya que involucran una gran cantidad de procesamiento externo a la CPU. Principalmente hay dos factores que hacen de los algoritmos externos diferentes de los algoritmos que veníamos viendo. Primero, el costo de acceder a un objeto es mucho mayor que cualquier costo asociado a los cálculos en la CPU. En general, podemos pensar que el tiempo de acceso de un dispositivo más rápido a uno más lento hay una constante de 10^6 de veces más lento. En segundo lugar, incluso con su mayor costo, hay varias restricciones en el acceso, dependiendo del tipo de memoria externa usada: por ejemplo, los registros en una cinta magnética solo pueden ser accedidos de forma secuencial.

Es por este motivo que se buscan algoritmos y estructuras de datos que minimicen la cantidad de operaciones sobre las memorias secundarias, a pesar de que eso resulte en más operaciones en la memoria principal. La idea es que el análisis asintótico sigue valiendo, pero cuando las constantes son tan grandes, vale la pena pensar soluciones alternativas. En particular, nos vamos a centrar en el problema de ordenamiento y búsqueda, pero con la diferencia que el volumen de los datos hace que los mismos se encuentren, necesariamente, en la memoria secundaria, y lo que buscamos es minimizar la cantidad de veces que un elemento es movido de la memoria secundaria a la memoria principal, y que estas transferencias se realicen de la forma más eficiente que permita el hardware.

Vamos a centrarnos en los métodos básicos de ordenamiento sobre cintas magnéticas y discos, porque es probable que estos dispositivos sigan siendo de uso masivo y además ilustran dos modos fundamentalmente distintos de acceso que caracterizan muchos sistemas de memoria externa. Para muchos de los dispositivos de memoria secundaria, el tiempo de acceso está dominado por el tiempo que se tarda en posicionar para leer el dato, más aún considerando que los datos se pueden acceder de a bloques en forma simultánea. Por ejemplo, en un disco rígido se leen de a bloques (de tamaño fijo) de bytes, y no de a bytes sueltos. Por este motivo, conviene pensar en algoritmos que accedan a datos contiguos en la memoria secundaria.

6.1. Ordenamiento Externo basado en Sort-Merge

La mayoría de los métodos de ordenamiento externo usan la siguiente estrategia general, conocida como Sort-Merge (Ordenación-Fusión): primero hacer una pasada sobre el archivo a ser ordenado, dividiéndolo en bloques de tamaño similar al de la memoria interna, y ordenar estos bloques. Luego combinar los bloques ordenados, haciendo varias pasadas por el archivo, obteniendo bloques ordenados cada vez más largos hasta terminar con todo el archivo ordenado. Los datos suelen ser accedidos de forma secuencial,

lo que hace a este método apropiado para la mayoría de los dispositivos externos. Los algoritmos de ordenamiento se esfuerzan por reducir el número de pasadas sobre el archivo y reducir el costo de una sola pasada para que sea lo más cercano al costo de copia posible.

Como la mayoría de los costos en un método de ordenamiento externos se debe a la entrada / salida, podemos dar una medida aproximada del costo de un algoritmo de sort-merge a partir de contar el número de veces en las que cada palabra en el archivo es leída o escrita. Por lo tanto, nuestro objetivo va a ser minimizar el número de pasadas que hacemos sobre el archivo. Dentro de esta familia de algoritmos de ordenamiento externo, vamos a centrarnos en los siguientes algoritmos:

- Fusión Múltiple Equilibrada (Balanced Multiway Merging).
- Selección por sustitución.
- Fusión Polifásica.

6.1.1. Fusión Múltiple Equilibrada

Para empezar, vamos a seguir los distintos pasos del procedimiento más simple basado en sort-merge para un archivo de ejemplo. Supongamos que tenemos que ordenar el siguiente archivo:

”EJEMPLODEORDENACIONFUSION”

Los registros están en una cinta y solo pueden ser accedidos secuencialmente (el segundo registro no puede ser leído antes que el primero, etc). Además, en memoria solo hay espacio para un número fijo de registros (vamos a suponer 3 registros), pero disponemos de todas las cintas auxiliares que necesitamos.

El primer paso para el algoritmo de FME consiste en recorrer secuencialmente el archivo de entrada, leyendo 3 registros de forma simultánea, ordenándolos en memoria principal, para finalmente escribirlos ordenados en una cinta auxiliar. Ahora, para poder combinar estos bloques, estos deben estar en cintas diferentes. Si queremos hacer una combinación de 3 vías, entonces usaríamos tres cintas auxiliares distintas (escribiendo el resultado del i -ésimo bloque leído en la cinta i módulo 3). Si el archivo no es múltiplo de 3 registros, simplemente el último bloque va a quedar con uno o dos elementos vacíos.

Primera Pasada

Entrada: Archivo a ordenar

```

1 while Haya elementos en la cinta de entrada do
2   for  $i = 1, \dots, 3$  do
3     Leer de a bloque 3 registros, ordenarlos en RAM, y escribirlos en la cinta  $i$ .

```

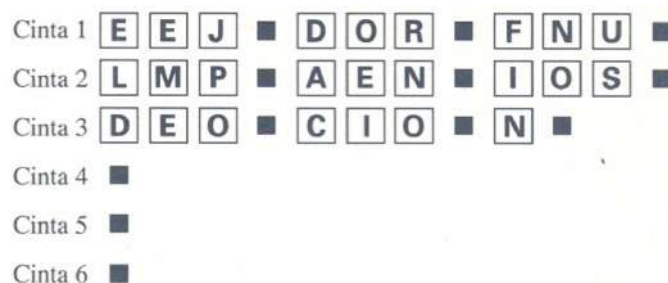


Figura 6.1: Fusión equilibrada de tres vías: resultado de la primera pasada.

Ahora estamos en condiciones de empezar a combinar los bloques ordenados de tamaño 3. Primero leemos el primer registro de cada cinta auxiliar y escribimos como salida aquella con la clave más chica entre los tres registros. Si hubiésemos usado más de 3 cintas para almacenar los registros, no podríamos hacer este paso al no alcanzarnos la capacidad de la memoria principal para hacer la fusión.

Luego se lee el siguiente registro de la misma cinta del registro que acaba de salir, para volver a escribir la clave más chica. Cuando llegamos al final de un bloque de tres palabras, entonces esa cinta es ignorada hasta que los bloques de las otras dos cintas se terminen de procesar, y los nueve registros se hayan escrito. El proceso se repite para combinar los segundos bloques de tres palabras de cada cinta en un bloque de nueve palabras (que se escribe sobre una cinta diferente, para estar lista para la siguiente combinación). Continuando, tenemos tres bloques largos ordenados que pueden combinarse para completar el ordenamiento.

Si tuviésemos un archivo mucho más largo con muchos bloques de tamaño 9 en cada cinta, entonces terminaríamos la segunda pasada con bloques ordenados de 27 elementos en las cintas 1, 2 y 3. Luego, una tercera pasada produciría bloques de tamaño 81 en las cintas 4, 5 y 6, y así siguiendo. En conclusión, necesitamos seis cintas para ordenar un archivo de tamaño arbitrario: tres para usarlas como entrada y tres para usarlas como salida de cada una de las combinaciones de 3 vías. Si no tenemos tantas cintas auxiliares, es posible utilizar tan solo 4 cintas: la salida se guardaría en una sola cinta, y los bloques de esta cinta se distribuiría en las tres cintas de entrada en cada paso de combinación.

Este método se conoce como *multiway merge*: es un algoritmo razonable para el ordenamiento externo y un buen punto inicial para la implementación de un ordenamiento externo. Los algoritmos más sofisticados que veremos más adelante pueden ordenar el doble de rápido, pero no mucho más (aunque cuando tenemos tiempos de ejecución muy altos, incluso una pequeña mejora puede ser significativa).

Segunda Pasada

Entrada: Archivo a ordenar

```

1 while Haya bloques de 3 en las cintas 1, 2, 3 do
2   for  $i = 4, \dots, 6$  do
3     Fusionar los 3 bloques de longitud  $3^1$  armando uno de  $3^2$ , y almacenarlo en la cinta  $i$ .
```

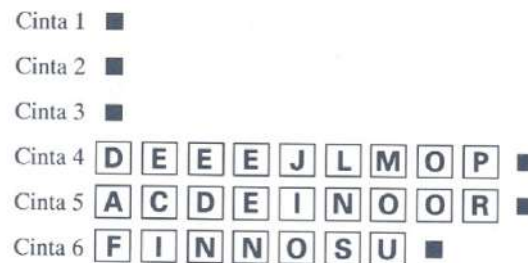


Figura 6.2: Fusión equilibrada de tres vías: resultado de la segunda pasada.

Supongamos que tenemos N palabras a ordenar y una memoria interna de tamaño M . Entonces, el primero paso de ordenamiento produce N/M bloques ordenados. Si hacemos una fusión de P vías en cada paso siguiente, entonces el número de pasadas sobre el archivo es alrededor de $\log_P(N/M)$, ya que cada pasada reduce el número de bloques ordenados en un factor P .

6.1.2. Selección por sustitución

Resulta que los detalles de implementación del método anterior pueden ser desarrollados de forma elegante y eficiente usando colas de prioridad. Primero, vamos a ver que las colas de prioridad proveen una forma natural de implementar una *multiway merge*. La idea va a ser utilizar una **cola de prioridad** de tamaño M para implementar ambas partes de la FME. La ventaja que nos da la utilización de una cola de prioridad es que vamos a poder trabajar con una mayor cantidad de bloques ordenados, superando el límite impuesto por la capacidad de la memoria principal. Notemos que el caso anterior, solo podíamos ordenar de a 3 bloques, al contar con una memoria de solo 3 registros.

Lo único que se necesita para realizar la operación de fusión a P vías es poder darle salida al

elemento más chico perteneciente a los bloques a fusionar. Una vez extraído ese elemento, se reemplaza con el siguiente elemento del bloque al que pertenecía. Entonces, la idea va a ser utilizar una cola de prioridad para armar esos bloques, y así vamos haciendo *pasar* a los elementos del archivo a través de la cola de prioridad, sustituyendo al elemento más chico por el siguiente del bloque. Cuando se termina de extraer todos los elementos de un bloque, se coloca un centinela en el heap, el cual es considerado mayor a cualquier otra clave. Cuando todos los elementos en el heap son centinelas, la fusión ha sido completada. Esta manera de usar colas de prioridad a veces se llama *replacement selection* (selección por sustitución).

Por lo tanto, para hacer una fusión de P vías, podemos usar selección por sustitución sobre una cola de prioridad de tamaño P para encontrar en la memoria principal cada elemento a escribir en costo $\Theta(\log P)$. Esta mejora en eficiencia no es de particular interés práctico, ya que una implementación de fuerza bruta puede encontrar cada elemento en $O(P)$, y P suele ser tan chico como para que este costo sea irrelevante comparado con el costo de escribir el elemento en la cinta. La importancia real en el método de selección por sustitución es la manera en la que puede ser usado en la primera etapa del proceso de sort-merge: para formar los bloques ordenados iniciales que proveen la base de las pasadas de fusión.

La idea es pasar el input (desordenado) a través de una larga cola de prioridad, siempre sacando el elemento más chico en la cola de prioridad, como hacíamos antes, y siempre reemplazarlo con el siguiente elemento de la cinta a la que pertenecía, con una salvedad adicional: si el nuevo elemento es más chico que el que acaba de salir, entonces, no es posible que haya sido parte del bloque ordenado actual. Por lo tanto, lo que podemos hacer es marcarlo como miembro del siguiente bloque y tratarlo como si fuera mayor que todos los otros elementos en el bloque actual. Un elemento marcado de la misma manera que el elemento raíz se considera que está en el bloque actual, y que el resto pertenece al siguiente bloque. De esta manera, evitamos el uso de centinelas, y permitimos generar salidas más largas. Cuando un elemento marcado llega a la cima de la cola de prioridad, el viejo bloque ya terminó y un nuevo bloque acaba de empezar.

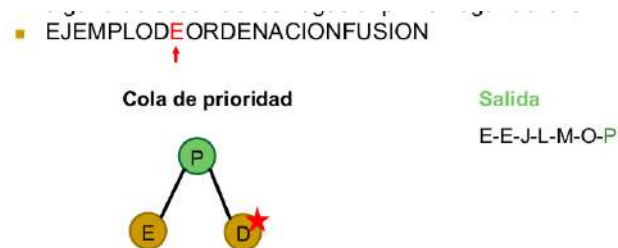


Figura 6.3: Ejemplo de Selección por sustitución.

Una vez tenemos las distintas cintas ya ordenadas, el algoritmo continúa aplicando fusiones como en el caso anterior. Se puede demostrar que, si las claves son aleatorias, las secuencias ordenadas creadas con este algoritmo tienen el doble de longitud que las que podrían producirse usando un método interno. El efecto práctico que ganamos con esto es que nos ahorramos una pasada sobre la secuencia: en lugar de iniciar con secuencias ordenadas de tamaño cercano al de la memoria principal, para luego hacer un paso de fusión para producir secuencias del doble del tamaño, podemos empezar con secuencias del doble del tamaño de la memoria, utilizando una cola de prioridad de tamaño M . Si existe algún orden entre las claves, las secuencias serán mucho más largas. Por ejemplo, si ninguna clave tiene más de M claves mayores antes de ella en el archivo, el archivo será completamente ordenado en una sola pasada, y no se necesitará de hacer fusiones. Esta es la razón práctica más importante por la cual usar este método.

En resumen, el método de selección por sustitución puede ser usado para ambos pasos de ordenar y fusionar de una Fusión Múltiple Balanceada. Para ordenar N registros usando una memoria interna de tamaño M y $P+1$ cintas auxiliares, primero usamos selección por sustitución con una cola de prioridad M para producir secuencias ordenadas iniciales de longitud $2M$ o más, y luego usar selección por sustitución con una cola de prioridad P para hacer $\log_p(N/2M)$ (o menos) pasadas de fusión.

6.1.3. Fusión Polifásica

Uno de los problemas con la fusión múltiple balanceada es que requiere de una cantidad excesiva de cintas auxiliares (o hacer muchas más copias de las necesarias). Para una fusión de P vías debemos usar $2P$ cintas (P como entrada y P como salida) o debemos copiar casi todo el archivo de una sola cinta de salida para P cintas de entrada entre cada pasada de fusión, que termina duplicando la cantidad de pasadas. Varios algoritmos de ordenamiento para cintas se han inventado para eliminar virtualmente todas estas copias a través de cambiar la manera en la que los bloques más cortos son ordenados y fusionados. El método más conocido es la *Fusión Polifásica* (polyphase merging).

La idea básica detrás de la fusión polifásica es distribuir los bloques ordenados producidos por selección por sustitución de manera algo dispareja entre las cintas disponibles, dejando una completamente vacía, y luego aplicar una estrategia de "fusión hasta el vaciado", momento en el que una de las cintas de salida y una de entrada cambian de roles (se rebobinan las cintas). Una de las ventajas que tiene este algoritmo es que nos permite trabajar con un número arbitrario de cintas que dispongamos para el ordenamiento.

Por ejemplo, supongamos que tenemos tan solo tres cintas, y empezamos con la siguiente configuración inicial de bloques ordenados en cintas. Primero aplicamos selección por sustitución sobre el archivo con una memoria interna que puede almacenar tan solo dos registros. Luego de tres fusiones a dos vías de las cintas 1 y 2 a la cinta 3, la segunda cinta termina vacía y nos quedamos con la siguiente configuración. Después, continuamos con dos fusiones de 2 vías entre las cintas 1 y 3 con salida a la cinta 2, quedando la primera cinta vacía.

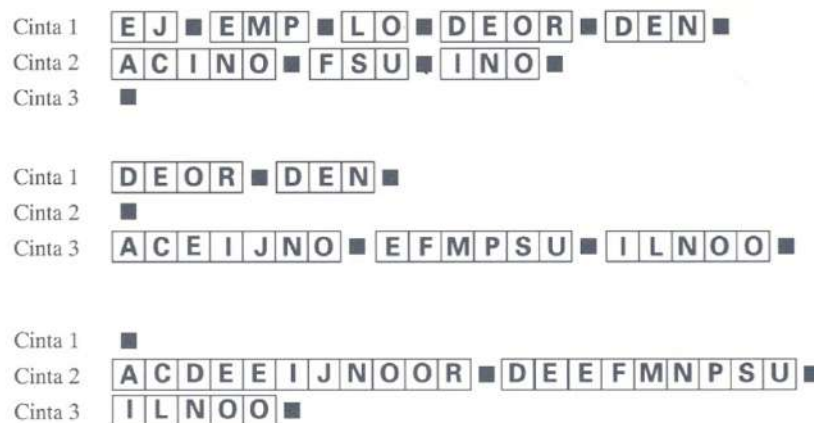


Figura 6.4: Etapas iniciales de una fusión polifásica con tres cintas.

El ordenamiento es completado en dos pasos más. Primero, una fusión a 2 vías entre las cintas 2 y 3 con salida a la cinta 1. Por último, hacemos una fusión a 2 vías entre las cintas 1 y 2 con salida a la cinta 3, terminando con todo el archivo ordenado en la cinta 3.

Esta estrategia de "fusión hasta el vaciado" puede ser extendida para funcionar con un número arbitrario de cintas. En general, la fusión polifásica solo es mejor que la fusión múltiple balanceada cuando P es chico. Para $P > 8$, es más probable que la fusión balanceada sea más rápida que la polifásica, y el efecto de usar la fusión polifásica cuando P es pequeño básicamente es ahorrarse dos cintas, ya que una fusión balanceada con 2 cintas extras sería más eficiente.

6.2. Búsqueda Externa

Los algoritmos de búsqueda apropiados para el acceso de objetos de archivos muy grandes tienen una inmensa importancia práctica. La búsqueda es una operación fundamental en grandes archivos, y ciertamente consume una fracción bastante significativa de los recursos de usados en muchas computadoras.

Vamos a centrarnos en los métodos de búsqueda en grandes discos de archivos, ya que las búsquedas sobre discos son de mayor interés práctico. Con los dispositivos secuenciales como las cintas, las búsquedas deben ser secuenciales, por lo que no podemos hacer mucho más allá del método trivial de montar la cinta y leerla hasta encontrar el objeto. Notablemente, los métodos que vamos a estudiar son capaces de encontrar un objeto en un disco de miles de millones de palabras en tan solo dos o tres accesos al disco.

Para muchas aplicaciones nos gustaría poder cambiar, agregar, borrar o acceder poca información dentro de archivos muy grandes. Los métodos que vamos a discutir tienen importancia práctica en la implementación de grandes *file systems* en los que cada archivo tiene un identificador único, y que tienen el propósito de dar soporte a acceso, inserción y borrado basado en ese identificador. Vamos a considerar que el disco está dividido en páginas, bloques contiguos de información que pueden ser accedidos de forma eficiente por el hardware del disco. Cada página va a contener muchos registros; nuestra tarea es organizar los registros dentro de las páginas de manera tal que cualquier registro pueda ser accedido tan solo leyendo unas pocas páginas. Vamos a asumir que el tiempo de E/S requerido para leer una página completa domina al tiempo de procesamiento requerido para hacer cualquier cómputo que involucre a esa página. La idea, nuevamente, va a ser saber aprovechar todos los datos que se leen de a bloques, y así minimizar la cantidad de accesos a memoria secundaria. Los métodos que vamos a ver son:

- Acceso Secuencial Indexado (ISAM): solución histórica.
- Árboles B: sirven tanto para memoria búsqueda en memoria principal como para búsqueda en memoria secundaria.
- Hashing Extensible.

6.2.1. Acceso Secuencial Indexado

La búsqueda secuencial consiste en tener a los registros almacenados en orden creciente según sus claves, y las búsquedas se realizan simplemente leyendo los registros uno atrás del otro hasta encontrar una clave que sea mayor o igual que la clave buscada. Para mejorar considerablemente la velocidad de la búsqueda, podemos mantener un "índice", que se almacenará como primer página de cada disco, que nos diga qué claves pertenecen a las páginas de ese disco, permitiendo la posibilidad de mirar únicamente este índice para saber si está presente una clave, sin tener que acceder a todas las páginas del disco. Por cada página, tenemos en el índice una entrada que nos indica cuál es la mayor clave que se encuentra en la página anterior (la primer entrada se usa para saber cuál es la mayor clave de la última página del disco anterior). De este modo, podemos saber si una clave está o no en un conjunto de discos a partir de hacer un acceso por disco, obteniendo su índice y un acceso adicional para obtener dicha página.

Estos índices son acompañados de un "índice maestro", que es una estructura que nos dice para cada disco, cuál es la última clave que contiene. Por ejemplo, el índice maestro nos diría que el disco 1 contiene claves $\leq E$, el disco 2 claves $\leq O$ (pero $\geq E$), etc. Este suele ser lo suficientemente chico como para que fácilmente se guarde en la memoria principal, de manera tal que la mayoría de los registros puedan ser encontrados en tan solo dos accesos de página: uno para el índice de algún disco apropiado y uno para la página que contiene al registro. Notemos que estamos realizando una búsqueda secuencial sobre el índice maestro, por lo que una posible mejora sería mantener ordenados los discos de manera tal de que podamos realizar búsqueda binaria sobre este, aunque esto no es importante, al tratarse de accesos a memoria principal.

Como se combina una organización secuencial de claves con el acceso indexado, esta organización se llama *Acceso Secuencial Indexado*. Este método es el apropiado para aplicaciones en las que los cambios sobre la base de datos son poco frecuentes, ya que incluso agregar una sola clave puede implicar una reorganización completa de los discos. Sin embargo, nos podría ser de utilidad en caso de necesitar estructuras estáticas.

6.2.2. Árboles B

Una mejor manera de manejar la búsqueda en una situación dinámica es usar árboles balanceados (que se almacenará en la memoria secundaria). Consideremos un árbol binario de búsqueda grande, e imaginemos que se ha almacenado en un disco. Si buscamos sobre este árbol de la manera en la que

veníamos viendo para algoritmos internos, tendríamos que hacer cerca de $\log n$ accesos a disco antes de completar la búsqueda. Ahora, supongamos que dividimos el archivo en páginas 7-nodo; si accedemos una página a la vez, entonces necesitaríamos tan solo un tercio de los accesos que teníamos antes, por lo que la búsqueda sería tres veces más rápida. Este tipo de árboles son llamado árboles B.

Agrupar los nodos en páginas es el cambio esencial entre un árbol binario y un árbol B, con 8 vías de ramificación en cada página-nodo. Si tenemos páginas aún más grandes, con una ramificación de 128 vías luego de cada acceso, podríamos encontrar cualquier clave en una tabla de un millón de entradas en tan solo tres accesos a disco. Podemos mantener la página raíz en todo momento en la memoria principal, de manera tal que reduzcamos la cantidad de accesos a 2 (la memoria principal tendría que ser capaz de almacenar dos nodos-páginas).

Los árboles B son árboles perfectamente balanceados, es decir, todas las hojas están a la misma distancia de la raíz. La diferencia con los árboles que habíamos visto anteriormente, es que no son árboles estrictamente binarios, sino que van cambiando la aridad de las ramas. El objetivo es que solo se necesite acceder a una única página de disco por cada nivel del árbol. Para que el árbol no sea demasiado alto, se necesita que haya mucha ramificación, por lo que se almacenan muchas claves en cada nodo.

Los Árboles 2-3-4 son un caso particular de los Árboles B, en donde los nodos pueden contener 1, 2, o 3 claves. Un i -nodo contiene $i - 1$ claves, y de este se van a desprender i subárboles. La idea es que si en un nodo se tienen dos claves $c_1 < c_2$, el primer subárbol contenga claves menores que c_1 , el segundo subárbol contenga claves entre c_1 y c_2 , y el tercer subárbol contenga claves mayores que c_2 . Decimos que el nodo está **saturado** si no podemos agregarle más claves.

Cuando hacemos una inserción, tenemos que preservar el invariante de los árboles 2-3-4. Para ello, lo que se hace es buscar en el árbol la posición que le correspondería a la nueva clave. Si el nodo es un 2-nodo (o un 3-nodo), simplemente agregamos la clave al nodo, quedándonos con un 3-nodo (o un 4-nodo). En caso de que se trate de un 4-nodo, no podemos agregar la clave al nodo. Por lo tanto, lo que vamos a hacer es dividir el 4-nodo en dos 2-nodos, y vamos a pasar una de sus claves hacia arriba en el árbol. La idea es que un 4-nodo es muy parecido a un árbol binario de tres nodos.

Supongamos que tenemos que insertar las claves A S E A R C H I N G E X A M P L E en un árbol inicialmente vacío (en ese orden). Empezamos con un 2-nodo (A), luego un 3-nodo (A S), luego un 4-nodo (A E S). Cuando tenemos que poner una segunda A en el 4-nodo, el algoritmo divide el 4-nodo formando un árbol binario antes de insertar la clave A, de manera tal que ahora haya un espacio para A en el fondo del árbol.

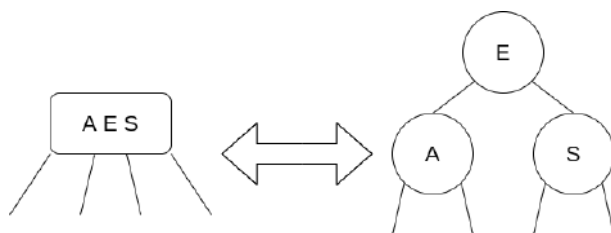


Figura 6.5: Equivalencia entre 4-nodo y árbol binario

En definitiva, cuando tenemos un 4-nodo, este se divide en dos 2-nodos para hacer lugar a la nueva clave, reubicando la clave del medio al nivel de arriba. El problema que nos falta resolver es qué hacemos cuando tenemos que dividir un 4-nodo cuyo padre es también un 4-nodo. Un método sería repetir el procedimiento sobre el padre de manera recursiva, hasta llegar a algún nivel que no tenga su nodo saturado o hasta llegar a la raíz (y simplemente agregar un nuevo nivel como vimos para el caso de A E S). Este tipo de árboles 2-3-4 se conocen como *bottom-up 2-3-4 trees*.

Otra posibilidad es asegurarse que el padre de cualquier nodo que veamos en el camino de bajada no sea un 4-nodo. Si alguno de ellos fuera un 4-nodo, lo dividimos de forma preventiva. Específicamente, cada vez que nos encontremos con un 2-nodo conectado a un 4-nodo, lo deberíamos transformar en un 3-nodo conectado a dos 2-nodos; y cada vez que nos encontremos en el camino de bajada a un 3-nodo conectado

a un 4-nodo, deberíamos transformarlos en un 4-nodo conectado a dos 2-nodos. Esta alternativa funciona mejor en la práctica, principalmente para accesos concurrentes a los datos¹.

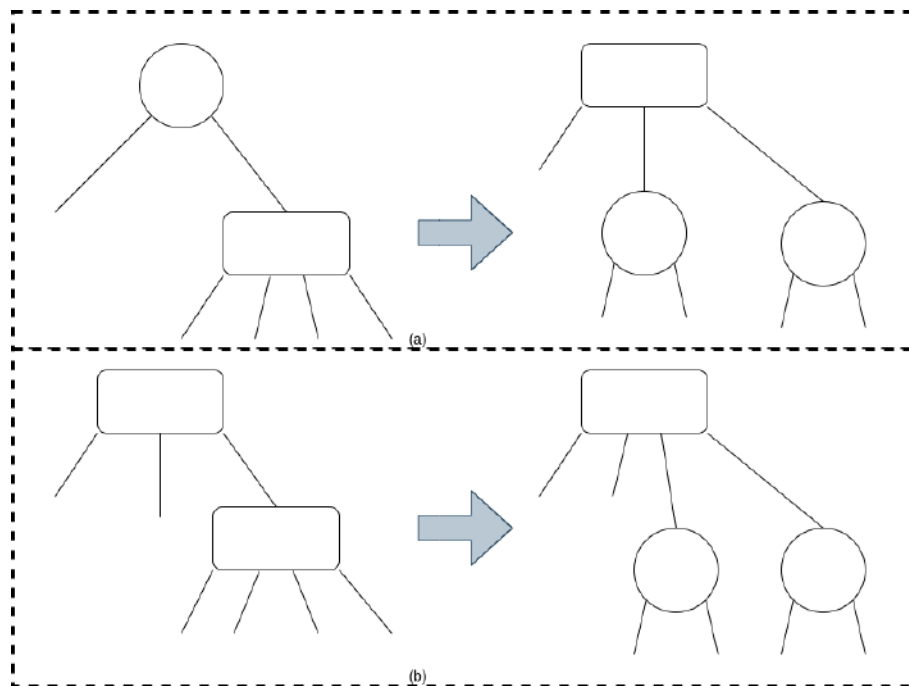


Figura 6.6: (a) 2-nodo conectado a 4-nodo. (b) 3-nodo conectado a 4-nodo.

Estas transformaciones son completamente locales: ninguna parte del árbol debe ser examinada o modificada salvo por aquellas que vimos. Cada una de las transformaciones suben una de las claves de un 4-nodo a su padre en el árbol, reestructurando los enlaces de forma acorde. Notemos que no tenemos que preocuparnos explícitamente acerca de si el padre es un 4-nodo, ya que nuestras transformaciones aseguran que a medida que pasamos por cada nodo en el árbol, terminamos con que ninguno es un 4-nodo. En particular, cuando salimos del último nivel del árbol, no estamos en un 4-nodo, y podemos insertar directamente el nuevo nodo ya sea transformando un 2-nodo a un 3-nodo o un 3-nodo a un 4-nodo. De hecho, es conveniente tratar la inserción como una división de un 4-nodo imaginario en el último nivel del árbol, que termina haciendo subir a la nueva clave a insertar. Cuando la raíz se vuelve un 4-nodo, vamos a dividirlo en tres 2-nodos, como hicimos en el primer caso A E S, haciendo que el árbol crezca un nivel para arriba.

Este algoritmo nos da una manera de hacer búsquedas e inserciones en los árboles 2-3-4. Como los 4-nodos se dividen en el camino de bajada, estos árboles se llaman *top-down* 2-3-4 trees. Lo interesante es que, incluso sin habernos preocupado por el balanceo, el árbol resultante está perfectamente balanceado. La distancia de la raíz a cualquier nodo hoja es siempre la misma, lo que implica que el tiempo requerido para una búsqueda o una inserción es siempre proporcional a $\log n$. Para ver que el árbol siempre queda perfectamente balanceado simplemente tenemos que darnos cuenta que las operaciones que realizamos no tienen ningún efecto sobre la distancia entre ninguna hoja a la raíz, exceptuando el caso en el que dividimos la raíz, y esto causa que la distancia de todos los nodos a la raíz aumente en uno.

Los algoritmos top-down que vimos para los árboles 2-3-4 se extienden fácilmente para manejar nodos con más claves, con la diferencia de que cada nodo puede tener hasta $M - 1$ claves y M subárboles. La búsqueda procede de una manera análoga a la de los árboles 2-3-4: para movernos de un nodo al siguiente, primero encontramos el intervalo adecuado para la clave en el nodo actual y luego salimos por el eje correspondiente. Continuamos de esta manera hasta que llegamos al nodo hoja correspondiente, para luego insertar la nueva clave en el último nodo interno alcanzado. Al igual que en los árboles top-down

¹No encontré ninguna explicación sobre esto. Al rebalancear de forma top-down, liberamos de cualquier race condition a $k - 1$ subárboles de en el primer paso, siendo k la aridad de la raíz. En cambio, si rebalanceamos de forma bottom-up, en el primer paso tan solo nos libramos de race conditions sobre el subárbol de altura 2 en el que estamos.

2-3-4, es necesario dividir los nodos que están saturados en el camino de bajada del árbol: en cualquier momento en el que veamos un k -nodo unido a un M nodo, lo reemplazamos por un $k + 1$ -nodo unido a dos $M/2$ nodos. Esto garantiza que cuando se llega al piso del árbol haya espacio para insertar un nuevo nodo.

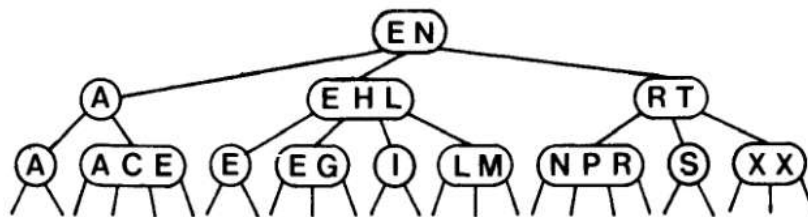


Figura 6.7: Ejemplo de árbol-B construido para $M = 4$ y el archivo E X T E R N A L S E A R C H I N G E X A M P L E

Este árbol tiene 13 nodos, cada uno correspondiente a una página del disco (la asignación de nodos a páginas de discos requiere de estrategias sofisticadas que no vamos a ver). Recordemos que los discos tenían capacidad de almacenar tres páginas de cuatro registros cada una. La elección de $M = 4$ es para enfatizar el siguiente punto: en las páginas debemos almacenar tanto los registros como los enlaces a los siguientes nodo. La cantidad total de espacio utilizado va a depender del tamaño relativo entre los registros y los enlaces.

Muchos nodos en el último nivel del árbol B descrito contienen muchos enlaces sin usar que podrían ser eliminados marcando a estos nodos hoja de alguna manera. Es posible optimizar esta estructura utilizando valores de M mucho más altos para los nodos internos del árbol si en lugar de almacenar los registros, tan solo almacenaran las claves, al igual que hacíamos en el método de acceso secuencial indexado (almacenando todos los registros en los nodos hoja). Por ejemplo, suponiendo que podemos colocar hasta 7 claves y 8 enlaces en una página, podríamos usar $M = 8$ para los nodos internos y $M = 5$ para los nodos hoja. Esta idea nos dejaría con el siguiente árbol:

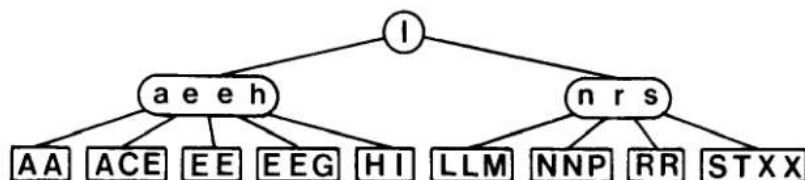


Figura 6.8: Optimización de árbol B

El efecto en una situación real sería mucho más dramático, ya que el factor de ramificación del árbol se incrementa aproximadamente según la relación entre el tamaño de un registro y el tamaño de una clave, que suele ser grande. Además, con esta organización, el "índice" que contiene claves y enlaces puede ser separado de los registros, como pasaba en el método de acceso secuencial indexado. De esta manera, tendríamos dos valores de M : uno para los nodos internos que determina el factor de ramificación del árbol (M_I) y uno para los nodos hoja que determina la posición de cada registro en el disco (M_B). M_B suele ser la cantidad de registros que entran en una página, ya que el objetivo de la búsqueda era encontrar la página que contenía al registro deseado. M_I se suele elegir lo suficientemente grande como para que la altura del árbol B sea de tan solo 3 niveles. Normalmente, M_B se corresponde con el número de claves que entran en 2 o 4 páginas consecutivas. Esta variante de los árboles B se conoce como árboles B^+ . Recordemos que como todas las búsquedas involucran el acceso al nodo raíz, este se suele mantener en la memoria principal, de manera tal que podamos encontrar cualquier elemento en el archivo en tan solo 2 accesos al disco.

Existen otras optimizaciones que se pueden hacer a los árboles B. Por ejemplo, una variación importante bajo el contexto de búsqueda externa es la idea de *overflow*. La idea es mejorar el algoritmo de inserción a través de resistir la tentación de dividir a los nodos de forma tan frecuente; para evitar la división, se hace una rotación local en la que le pasamos claves a un nodo hermano. Esto lleva a una mejor utilización del espacio en los nodos (asegura que se usan $2/3$ de la capacidad de cada nodo, en lugar de $1/2$), que suele ser de gran importancia cuando trabajamos con aplicaciones de búsqueda sobre discos de gran escala (aunque ralentiza las inserciones). Esta variante se conoce como árboles B^* .

6.2.3. Hashing Extensible

Una alternativa a los árboles B es el método de hashing extensible. Este método garantiza que cualquier búsqueda no requiera más de dos accesos a disco. Al igual que en los árboles B, los registros son almacenados en páginas que cuando se llenan, se dividen en dos partes (potencialmente utilizando el mismo algoritmo que vimos para Árboles B); al igual que en acceso secuencial indexado, vamos a mantener un índice en memoria principal que nos permite encontrar la página que contiene el registro asociado a nuestra clave de búsqueda.

Para ver cómo funciona el método de hashing extensible, vamos a considerar cómo es que maneja las inserciones sucesivas del siguiente ejemplo: E X T E R N A L S E A R C H I N G E X A M P L E, usando páginas con capacidad de cuatro registros y dos discos. Empezamos con un "índice" de una sola entrada, un puntero a la página que almacenará los registros. Los primeros cuatro registros entran en la página, por lo que nos queda la siguiente estructura:

- Disco 1: 20.
- Disco 2: EETX.

El directorio en el disco 1 nos dice que todos los registros están en la página 0 del disco 2, donde se almacenan en forma ordenada según sus claves (como las claves se ordenan en memoria principal, no resulta costoso). Para decidir a dónde va cada clave, se utiliza una sucesión de funciones hash h_1, h_2, \dots tales que h_i es una función hash binaria, que dado una clave cualquiera devuelve 0 o 1. En este caso, E = 00101², T = 10100, X = 11000. Ahora, la página está llena, y debe dividirse para poder agregar la clave R = 10010. La estrategia es simple: ponemos los registros con tales que $h_1(k) = 0$ en una página y en la otra ponemos los registros con $h_1(k) = 1$. Esto requiere que dupliquemos el tamaño del directorio, y que movamos la mitad de las claves de la página 0 del disco a una nueva página, dejando la siguiente estructura:

- Disco 1: 20 21
- Disco 2: EE RTX

Ahora el directorio en el disco 1 nos dice que todos los registros están repartidos entre la página 0 del disco 2 y la página 1 del disco 2. Ahora podemos agregar N = 01110 y A = 00001, quedándonos:

- Disco 1: 20 21
- Disco 2: AEEN RTX

Al momento de agregar L = 01100, otra división es necesaria. Prosiguiendo de la misma manera que en la primera división, hacemos espacio para L = 01100 separando en dos la primera página, una para las claves que empiezan con 00 y otra para las claves que empiezan con 01. Además, duplicamos el tamaño del directorio, quedándonos la siguiente estructura:

- Disco 1: 20 21 22 22
- Disco 2: AEE LN RTX

De esta manera, primero debemos acceder al directorio usando los primeros dos bits de la clave (los que empiezan con 00 están en el disco 2 página 0, los que empiezan con 01 están en el disco 2 página 1, los

²O sea, $h_1(E) = 0, h_2(E) = 0, h_3(E) = 1, h_4(E) = 0, h_5(E) = 1$.

que empiezan con 10 están en el disco 2 página 2, los que empiezan con 11 están en el disco 2 página 2), y una vez sabemos en qué página están, podemos acceder a la página que contiene al registro en tan solo 2 accesos a disco.

En general, la estructura construida por hashing extensible consiste en un **directorio** de 2^d palabras (siendo cada una un patrón de d bits) y un conjunto de hojas-página que contienen todos los registros con claves que empiezan con cierto patrón de bits (de longitud menor o igual a d). Una búsqueda involucra el uso de los primeros d bits de la clave para indexar el directorio, el cual apunta a las hojas-página. Luego, la hoja-página referenciada es accedida y se busca (usando cualquier estrategia) para encontrar el registro en sí. Una hoja-página puede ser apuntada por más de una entrada de directorio (como vimos en el ejemplo con las entradas 10 y 11).

El directorio solo contiene punteros a páginas. Estos suelen ser más pequeños que los registros, por lo que vamos a tener muchas más entradas en cada página que registros. Cuando el directorio no cabe en una sola página, mantenemos un nodo raíz en memoria que nos indica dónde están las páginas del directorio, usando el mismo esquema de indexado. Por ejemplo, si el directorio ocupa dos páginas, el nodo raíz podría contener dos entradas "10 11", indicando que el directorio para todas los registros con claves que empiezan con 0 están en la página 0 del disco 1, y que el directorio para todos los registros cuyas claves empiezan con 1 están en la página 1 del disco 1.

En general, la inserción en una estructura de hashing extensible puede involucrar una de tres operaciones, luego de que se acceda a la página que podría contener la clave:

1. Si hay espacio en la página, el nuevo registro simplemente es insertado.
2. Si no hay espacio en la página, esta se debe dividir en dos.
3. Si el directorio tiene más de una entrada apuntando a esa página, una de las entradas puede reutilizarse para apuntar a la nueva página. Sino, es necesario duplicar el tamaño del directorio.

Si por causa de un borrado la unión de dos páginas hermanas entra en una única página, lo que se hace es unirlas en una sola página. En ambos casos, se requiere de a lo sumo 3 accesos a memoria secundaria. En la práctica, para evitar posibles secuencias de borrados e inserciones que terminen en divisiones y uniones sucesivas, lo que se hace es considerar que una página se llenó cuando se llega al 80% de capacidad, y se considera que dos páginas hermanas se deben unir cuando la unión ocupe el 40% de la capacidad de una página. Asumiendo que las funciones hash h_1, h_2, \dots distribuyen bien las claves (que todas las claves tengan la misma probabilidad de que salga 1 o 0), entonces la mitad de las claves van a estar a la derecha y la otra mitad a la izquierda, lo cual nos deja a la altura del árbol en $O(\log n)$ en tiempo promedio, siendo n la cantidad de entradas del directorio. Nuevamente, recordemos que este último análisis corresponde a la memoria principal y que solo necesitamos de una cantidad constante de accesos a memoria secundaria (a lo sumo 2).

Desde el punto de vista de hashing, podemos pensar que el algoritmo divide a los nodos para manejar colisiones: es por este motivo que el método se llama hashing extensible. Este método presenta una alternativa a los árboles B y al acceso secuencial indexado ya que usa exactamente dos accesos a disco para cada búsqueda (al igual en el acceso indexado), mientras que retiene la capacidad de inserciones eficientes (al igual que los árboles B).

Capítulo 7

Alternativas a los Árboles Balanceados

En este apartado, nos vamos a centrar en estudiar otras implementaciones eficientes de diccionarios, que no se basan en mantener balanceado el árbol. La primera estructura tienen la particularidad de que forma parte de los **algoritmos probabilísticos**, llamada Skip List. Este tipo de algoritmos renuncian al *determinismo*, es decir, a la propiedad de que siempre que le demos la misma entrada, se obtiene el mismo resultado. Al incorporar el azar a los algoritmos, perdemos en predecibilidad, pero ganamos en eficiencia y elegancia.

Dentro de la familia de los algoritmos probabilísticos, podemos diferenciar a los algoritmos del tipo Las Vegas: dan siempre el resultado correcto, pero su tiempo de ejecución $T(n)$ es una variable aleatoria. Para realizar un análisis de complejidad en el caso promedio, lo que hacemos es ver cuánto nos da la esperanza del tiempo de ejecución, es decir, $E(T(n)) \in O(f(n))$. Notemos que una de las versiones que vimos de Quick Sort es un algoritmo Las Vegas, ya que el pivote se elige al azar y se asegura que al finalizar la ejecución, el arreglo termina ordenado. Otro tipo de algoritmos probabilísticos son los del tipo Montecarlo: los resultados tienen cierta probabilidad de ser falsos, pero siempre tardan lo mismo.

La segunda estructura que vamos a ver tiene la particularidad de que no nos da una cota sobre el costo de una operación en el peor caso, sino que nos da una cota sobre el costo de una secuencia de operaciones. A este tipo de análisis, lo llamamos análisis del costo amortizado. Esta estructura son los Splay Trees, y nos dan una mayor eficiencia en el caso de que los accesos a las claves no sea uniforme, es decir, que haya claves que sean más frecuentemente accedidas que otras.

7.1. Skip Lists

Veamos ahora una estructura de datos randomizada llamada *skip lists*. Utiliza listas enlazadas ordenadas para implementar diccionarios. El problema que teníamos con las listas enlazadas era que, sin importar que la lista esté o no ordenada, no podíamos escapar la necesidad de recorrer *secuencialmente* la lista, por lo que las operaciones de búsqueda, inserción y borrado nos quedaba en $O(n)$ en el peor caso.

La idea de las Skip Lists es tener una forma de avanzar más rápido sobre la lista, sin tener que recorrer secuencialmente toda la lista. Para ello, lo que se hace es ir agregando nuevos *niveles* de punteros que apuntan a elementos cada vez más espaciados (exponencialmente). Con esto en mente, lo que vamos a hacer es darle al $c \cdot 2^i$ -ésimo nodo (con c distinto de toda potencia de 2) un puntero a los nodos en $2^0, 2^1, \dots, 2^i$ posiciones más adelante en la lista.

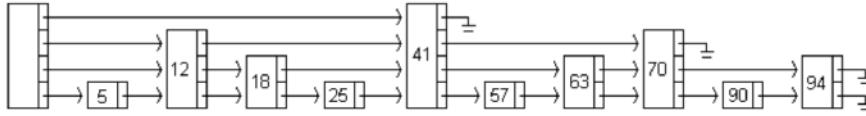


Figura 7.1: Ejemplo de Skip List.

De alguna manera, lo que nos estaría quedando son $\log n$ listas paralelas, cada una con una mayor distancia entre sus nodos sucesivos, de manera tal que por cada *nivel* que subimos en las listas, la distancia entre nodos aumenta de manera exponencial, y por tanto la cantidad de elementos en cada lista disminuye de manera exponencial. Para facilitar el acceso al i -ésimo nivel, lo que se hace es tener un primer nodo vacío, que tenga un puntero hacia el primer nodo de cada nivel. Si bien no estamos aumentando la cantidad de nodos de la lista, sí estamos aumentando la cantidad de punteros. Sin embargo, esto no es tan problemático, porque la cantidad total de punteros nos queda:

$$\begin{aligned}
 \lim_{k \rightarrow \infty} \sum_{i=0}^k \frac{n}{2^i} &= \lim_{k \rightarrow \infty} n \cdot \sum_{i=0}^k \left(\frac{1}{2}\right)^i \\
 &= \lim_{k \rightarrow \infty} n \cdot \frac{1 - \overbrace{\left(\frac{1}{2}\right)^{k+1}}^{\rightarrow 0}}{1 - \frac{1}{2}} \\
 &= \lim_{k \rightarrow \infty} n \cdot \frac{1}{0,5} = 2n
 \end{aligned}$$

donde n es la cantidad de elementos y k es la cantidad de niveles de la lista ($k = \log n$). Por lo tanto, solo necesitamos el doble de punteros que en una lista tradicional. Para realizar una búsqueda, a lo sumo se necesitan de 2 comparaciones por nivel, y como tenemos $\log n$ niveles, podemos concluir que el costo de la operación de búsqueda nos queda en $O(\log n)$. El problema principal con esta estructura es la rigidez que tiene frente a inserciones y borrados. Este tipo de Skip Lists se conocen como Skip Lists Perfectas, justamente por el hecho de que por cada nivel tenemos exactamente $2, 4, 8, \dots, 2^k$ punteros.

Para solucionar el problema de la rigidez, lo que se hace es someter a un proceso aleatorio la forma en la que se reconstruye la Skip List. Para ello, lo que hacemos es que al momento de insertar un nuevo elemento, lo ubicamos donde corresponda en la lista, y luego le asignamos a qué nivel pertenece de forma probabilística: tiramos una moneda de manera tal que si sale *cara*, promociona de nivel, y volvemos a iterar; si sale *crúz*, termina (todos los nodos empiezan siendo de nivel 1). En el caso promedio, sin tener que asumir propiedades probabilísticas sobre el input, intuitivamente nos va a quedar una Skip List con $O(\log n)$ niveles.

La ventaja con respecto a las estructuras de árboles AVL es que son mucho más simples de programar, extender y modificar. Por este motivo, suelen venir con mejores optimizaciones de implementación (por ejemplo, implementar el algoritmo de forma iterativa en lugar de forma recursiva). Además, las Skip Lists suelen tener un menor factor constante que las versiones poco optimizadas de los árboles AVL [8].

7.2. Splay Trees

Un splay tree es un árbol binario de búsqueda *auto-ajustante* en el que todas las operaciones básicas sobre árboles tienen un costo amortizado de $O(\log n)$, donde por costo amortizado nos referimos a que el tiempo por operación promedio sobre una secuencia de operaciones para el peor caso. Por lo tanto, los splay trees son tan eficientes como los árboles balanceados cuando es de interés medir el tiempo total de ejecución. Además, para secuencias de acceso lo suficientemente largas, los splay trees son tan eficientes como un árbol de búsqueda óptimo (su tiempo de ejecución es a lo sumo un múltiplo pequeño). La eficiencia de los splay trees no viene dado por una restricción explícita estructural, como sucede en los árboles balanceados, sino que su eficiencia se debe a la aplicación de una simple heurística de reestructuración, llamada *splaying*, cada vez que el árbol sea accedido. Extensiones de la operación de

splaying dan formas simplificadas de otras dos estructuras de datos: los árboles de búsqueda lexicográficos o multidimensionales y los link/cut trees[6].

La motivación de esta estructura es la observación de que los distintos tipos de árboles balanceados, si bien tienen una cota en el peor caso de $O(\log n)$, se pueden mejorar para el caso en el que el patrón de acceso no es uniforme, además del hecho de que necesitan de espacio adicional para almacenar la información de balanceo. Por ejemplo, cuando algunas claves son accedidas de forma mucho más frecuente que otras, nos gustaría que las importantes estén relativamente cerca de la raíz. Una estructura que nos garantiza esta propiedad son los árboles binario óptimos (ver sección 6.2.2 de [10]).

ABB óptimos: Recordemos que los árboles ABB cumplen con el invariante de que todo elemento en el subárbol derecho de un nodo es mayor este, y todo elemento en el subárbol izquierdo es menor. Si supiéramos la frecuencia de acceso a cada elemento, podríamos armar un árbol en el que los elementos más accedidos estén más cerca de la raíz. Esto se puede hacer en tiempo $O(n^2)$ a partir de un algoritmo basado en la técnica de Programación Dinámica [9]. Sin embargo, estar en una situación en la que sabemos la frecuencia de acceso de cada elemento y que esta no varíe a lo largo del tiempo lo vuelve un planteo poco realista.

Existen otras estructuras de datos, como los Biased search trees y los Finger search trees, que intentan mejorar sobre los ABB óptimos. Todas estas estructuras fueron diseñadas para reducir el costo en el peor caso de una operación. Sin embargo, en una aplicación típica de árboles de búsqueda, tenemos una secuencia de operaciones a realizar, y lo que importa es el tiempo total que tarda la secuencia en procesar, y no el tiempo individual de las operaciones. En este tipo de aplicaciones, un mejor objetivo sería reducir el tiempo amortizado de las operaciones.

Una forma de obtener eficiencia amortizada es usar una estructura de datos auto-ajustante. Vamos a permitir que la estructura esté en un estado arbitrario, pero en cada operación vamos a aplicar una regla simple de reestructuración que busca mejorar la eficiencia de operaciones futuras. En general, las ventajas que presenta los splay trees sobre los árboles balanceados es que resultan mucho más eficientes si el patrón de uso es sesgado (algunas claves se acceden mucho más que otras), utilizan menos espacio y los algoritmos son más sencillos. Sin embargo, esta estructura presenta dos posibles desventajas. En primer lugar, requieren de más ajustes locales, incluso durante los accesos, por lo que son poco eficientes para aplicaciones multi-threading. Además, las operaciones individuales podrían ser demasiado costosas, lo que sería un problema en aplicaciones real-time.

Frente al problema de falta de dinamismo en los ABB óptimos, los Splay Trees tratan de tender todo el tiempo al ABB óptimo en todo momento. Estos son más simples de implementar que los árboles AVL, al no tener la necesidad de verificar condiciones de balanceo y no requieren de memoria adicional (por ejemplo, para almacenar factores de balanceo). Decimos que son estructuras *auto-ajustantes*, porque se modifican así mismas cuando operamos sobre ellas. El ajuste que vamos a realizar es que cada vez que accedemos a un elemento, lo movemos hasta la raíz, mediante rotaciones parecidas a las del AVL, conocidas como *Splaying*. Además, si bien no se puede garantizar que cada operación cueste $O(\log n)$, sí se puede garantizar un costo de $O(m \cdot \log n)$ para m operaciones sucesivas, por lo que podemos decir que las operaciones de Splay Tree tienen costo **amortizado** de $O(\log n)$.

Para realizar un splay un árbol binario en un nodo interno x , empezamos en x y recorremos el camino hacia la raíz, en un orden que depende de la estructura del árbol. La operación de splay consiste en repetir el siguiente paso hasta que $p(x)$ quede indefinido, siendo $p(x)$ el padre de x y $p^2(x)$ su abuelo:

- (zig) Si x tiene un padre pero no un abuelo, rotamos en $p(x)$.
- (zig-zig) Si x tiene un abuelo y x y $p(x)$ son ambos hijos izquierdos o hijos derechos, rotamos en $p^2(x)$ y luego en $p(x)$.
- (zig-zag) Si x tiene un abuelo y x es un hijo izquierdo y $p(x)$ un hijo derecho, o viceversa, rotamos en $p(x)$ y luego en el nuevo padre de x (su antiguo abuelo).

En definitiva, el efecto de splaying es mover a x hasta la raíz del árbol mientras se acomoda al resto del camino original hasta x , permitiendo mantener a los elementos más accedidos lo más alto posible en el

árbol.

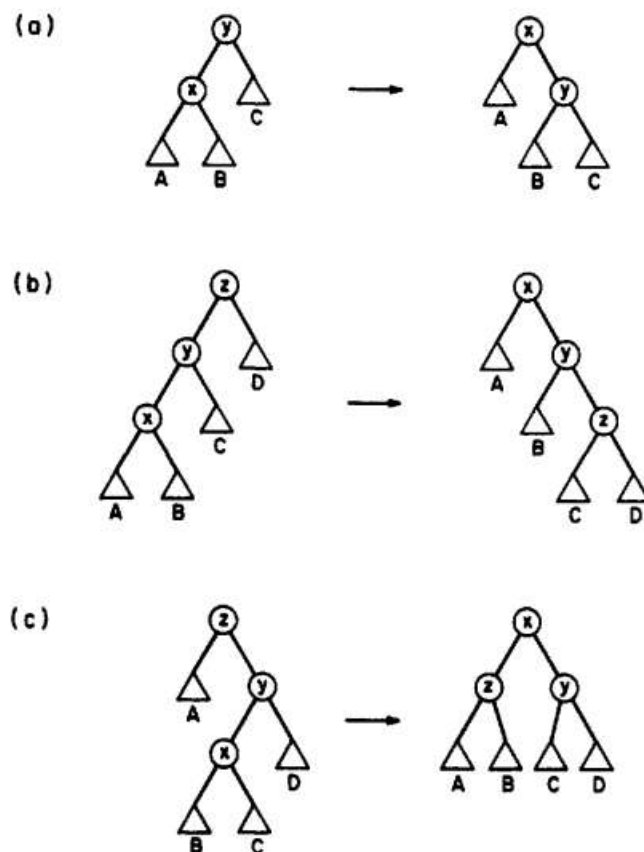


FIG. 4.9. *Cases of splay step. Each case has a symmetric variant (not shown). (a) Terminating single rotation. (b) Two single rotations. (c) Double rotation.*

Luego de acceder o insertar un elemento i , hacemos una operación de splay en i . Después de borrar un elemento i , hacemos un splay sobre su padre justo antes del borrado. En cada caso el tiempo de la operación es proporcional a la longitud del camino sobre el que se realiza el splay, al que llamamos *splay path*. Si bien hemos descrito a la operación de splay como bottom-up, existen variantes top-down, en donde efectuamos la operación de splay durante la bajada hasta i .

Recordemos que como efecto del splaying no solo se mueve el nodo que estamos accediendo (que lo movemos hacia la raíz), sino que todos los nodos del camino desde la raíz hasta el nodo accedido se mueven aproximadamente a la mitad de su profundidad anterior, a costa de que algunos *pocos* nodos bajen como máximo dos niveles en el árbol. El problema de que esta estructura se modifica incluso cuando la estamos consultando es que no sería práctica en aplicaciones multi-threading. Se sabe que el tiempo de acceso a los datos en un Splay Tree es a lo sumo un múltiplo pequeño del tiempo de acceso de los ABB óptimos estáticos, bajo un análisis de costo amortizado. Además, los Splay Trees se utilizan para resolver el problema de *cutting trees* de forma eficiente. Existen otras estructuras randomizadas que parecen ser más simples y rápidas que los splay trees, por ejemplo, los Cartesian trees o *treaps*.

Listas auto-ajustantes: Lo que queremos es tener una lista que mantenga más a mano a los elementos más recientemente accedidos. Esto se puede hacer aplicando una política Move-to-front (MTF), en la que se mueve el elemento accedido a la cabeza de la lista. Se puede demostrar que el costo total para una secuencia de m operaciones en una lista MTF es de a lo sumo el doble que el de cualquier implementación de diccionario usando listas.

Bibliografía

- [1] URL: <https://stackoverflow.com/questions/30713269/why-would-one-use-a-heap-over-a-self-balancing-binary-search-tree>.
- [2] *Advantages of Trie Data Structure*. 2021. URL: <https://www.geeksforgeeks.org/advantages-trie-data-structure/>.
- [3] Thomas H. Cormen et al. *Introduction to Algorithms*. The MIT Press. McGraw-Hill, 2001.
- [4] *Algorithms*. Addison Wesley, 1983.
- [5] Bratley P. Brassard G. *Fundamental of Algorithmics*. International series of monographs on physics. Prentice Hall, 1995.
- [6] R. E. Tarjan D. D. Sleator. «Self-Adjusting Binary Search Trees». En: *Commun. ACM* 32.3 (jul. de 1985), págs. 652-686. URL: <https://www.cs.cmu.edu/~sleator/papers/self-adjusting.pdf>.
- [7] *Data structures and network algorithms*. Society for Industrial y Applied Mathematics, 1987.
- [8] En: *Commun. ACM* 33.3 (1990). Ed. por Peter J. Denning, págs. 668-676. ISSN: 0001-0782. URL: <https://homepage.cs.uiowa.edu/~ghosh/skip.pdf>.
- [9] D. E. Knuth. «Optimum binary search trees». En: *Acta Informatica* 1.1 (mar. de 1971), págs. 14-25. ISSN: 1432-0525. DOI: [10.1007/BF00264289](https://doi.org/10.1007/BF00264289). URL: <https://doi.org/10.1007/BF00264289>.
- [10] D. E. Knuth. *The art of computer programming*. Vol. 3. Addison Wesley Professional, 1998.
- [11] J. D. Ullman V. Aho J. E. Hopcroft. *Data structures and algorithms*. Addison Wesley, 1985.