

CC4102/CC40A/CC53A - Diseno y Analisis de Algoritmos

Jeremy Barbay

18 May 2011

Índice

1. Tecnicas avanzadas de diseno y analisis de algoritmos (4 semanas = 8 charlas = 720mns)	1
1.1. Material relevante de los años previos:	1
1.2. Introduccion	1
1.3. Dominios discretos y finitos	2
1.3.1. Introduccion	2
1.3.2. Algoritmos de Ordenamiento (Counting Sort, Bucket sort, radix sort, string sort) . .	2
Counting Sort $O(\sigma + n)$	2
Bucket Sort $O(\sigma + n)$	2
Radix Sort $O(n \lg_n \sigma) = O(cn)$	3
MAYB Provechando de las repeticiones en el modelo de Comparaciones :MAYB: . .	3
MAYB String Sort :MAYB:	3
1.3.3. Busqueda por Interpolacion/Extrapolacion	3
1.3.4. Estructuras de Arboles con Dominio Discreto	4
Tries o Arboles Digitales	4
Arboles de Sufijos (y Arreglos de Sufijos)	5
1.3.5. Hashing	6
Introduccion	6
Paradoxe del cumpleaños (Probabilidad de Colision)	7
Hashing Abierto	7
Hashing Cerrado	8
Ideal Hashing	8
Universal Hashing	9
1.3.6. Hashing en memoria externa	9
1.4. Tecnicas de Analisis	10
1.4.1. Introduccion	10
1.4.2. Analisi amortizada	10
MATERIAL A LEER	10
Principio de Analisi amortizada	11
Analisi amortizada de Colas de Prioridades	12
Árbol biselado ("Splay Tree")	16
APUNTES	16
1.4.3. Analisi adaptativa	16
MATERIAL A LEER	16
Analisi en el peor caso: a dentro de que?	16
Busqueda Doblada: $1 + 2 \lceil \lg p \rceil$ comparaciones	16
Finger Search Tree: la busqueda doblada de los arboles de busqueda	17
APUNTES	17
Algoritmo de Ordenamiento adaptivos basicos	17
Merge Sort Adaptivo: Runs	17
Local Insertion Sort: Inv	17

	Another Insertion Sort: REM	17
	Computacion de la Union	17
	Computacion de la Interseccion	17
1.4.4.	Algoritmos en linea	17
	List Accessing	17
	Paginamiento Deterministico	20
	BONUS: "Ski Renting	22
1.4.5.	MAYB Complejidad Parametrizada	22
1.4.6.	PREGUNTAS [0/1] : PREGUNTAS:	22
	Analisis Amortizada: arreglo dinamico (Part 1)	22
	Analisis Amortizada: arreglo dinamico (Part 2)	23
	TODO Analisis Amortizada: arreglo dinamico (Part 3)	23
	Analisis Amortizada: arreglo dinamico (Part 4)	23
1.5.	Conclusion	23
1.6.	RESUMEN Unidad 3	23

1. Tecnicas avanzadas de disenio y analisis de algoritmos (4 semanas = 8 charlas = 720mns)

1.1. Material relevante de los años previos:

- Colas de Prioridades <http://www.leekillough.com/heaps/><http://www.leekillough.com/heaps/>
- Arboles 2-3 (para "Finger Search Trees"
 - <http://www.dcc.uchile.cl/bebustos/apuntes/cc3001/Diccionario/4><http://www.dcc.uchile.cl/bebustos/apuntes/cc3001/Diccionario/#4>
- <http://www.wimp.com/justcoincidence/><http://www.wimp.com/justcoincidence/> Birthday Paradox, in English.
- Interpolation Search
 - CC3001?
- Counting Sort
 - CLRS
 - CC3001

1.2. Introduccion

1. Dominios discretos y finitos

a) Busqueda en Dominios discretos y finitos

- Interpolacion/extrapolation
- Tries o arboles digitales
- Arboles y Arreglos de Sufijos
- Hash y Hash en memoria Secundaria

b) Algoritmos de Ordenamientos con universo finito

- Counting Sort
- Bucket Sort
- Radix Sort

2. Tecnicas de Analisis

- a) Analisis amortizada
 - tecnicas
 - colas de prioridades
 - splay arboles
- b) Analisis parametrizada
 - busqueda doblada y finger search trees
 - ordenamiento adaptivo
 - operaciones de conjuntos adaptivos
- c) Analisis de Algoritmos en linea
 - “ski renting” problema
 - analisis competitiva (“Competitive Analysis”)

1.3. Dominios discretos y finitos

1.3.1. Introduccion

- afuera del modelo de comparacion.

1.3.2. Algoritmos de Ordenamiento (Counting Sort, Bucket sort, radix sort, string sort)

Counting Sort $O(\sigma + n)$

1. for $j = 1$ to σ do $C[j] \leftarrow 0$
2. for $i = 1$ to n do $C[A[i]] ++$
3. $p \leftarrow 1$
4. for $j = 1$ to σ do
 - for $i = 1$ to $C[j]$ do
 - $A[p++] \leftarrow j$

Este algoritmo es bueno para ordenar multi conjuntos (donde cada elementos puede ser presente muchas veces), pero pobre para diccionarios, para cual es mejor usar la extension logica, Bucket Sort.

Bucket Sort $O(\sigma + n)$

1. for $j = 1$ to σ do $C[j] \leftarrow 0$
2. for $i = 1$ to n do $C[A[i]] ++$
3. $P[1] \leftarrow 1$
4. for $j \leftarrow 2$ to σ do
 - $P[j] \leftarrow P[j - 1] + C[j - 1]$
5. for $i \leftarrow 1$ to n
 - $B[P[A[i]]++] \leftarrow A[i]$

Este algoritmo es particularmente practica para ordenar llaves asociadas con objetos, donde dos llaves pueden ser asociadas con algunas valores distintas. Nota que el ordenamiento es **estable**.

Radix Sort $O(n \lg_n \sigma) = O(cn)$

- Considera un arreglo A de tamaño n sobre alfabeto σ
- si $\sigma = n$, se recuerdan que bucket sort puede ordenar A en $O(n)$
- si $\sigma = n^2$, bucket sort puede ordenar A en $O(n)$:
 - 1 vez con los $\lg n$ bits de la derecha
 - 1 vez con los $\lg n$ bits de la izquierda (utilizando la estabilidad de bucket sort)
- si $\sigma = n^c$, bucket sort puede ordenar A
 - en tiempo $O(cn)$
 - con espacio $O(n)$
El espacio se puede reducir a $2n + \sqrt{n}$ con $\lg n/2$ bits a cada iteracion de Bucketsort, cambiando la complejidad solamente por un factor de 2.
- En final, si A es de tamaño n sobre un alfabeto de tamaño σ , radix sort puede ordenar A en tiempo $O(n \lceil \frac{\lg \sigma}{\lg n} \rceil)$

MAYB Provechando de las repeticiones en el modelo de Comparaciones :MAYB:

- Se puede o no?
- Ordenar en nH_i comparaciones

MAYB String Sort :MAYB:

- Problema: Ordenar k strings sobre alfabeto $[\sigma]$, de largo total $n = \sum_i n_i$.
- Si $\sigma \leq k$, y cada string es de mismo tamaño.
 - Si utilizamos bucket-sort de la derecha a la izquierda,

podemos ordenar en tiempo $O(n)$, porque $O(n \lceil \lg \sigma / \lg l \rceil)$ y $\sigma < n$.

- Si $\sigma \in O(1)$
 - Radix Sort sobre c simbolos, donde c es el tamaño minima de una string, y iterar recursivamente sobre el restos de la strings mas grande con mismo prefijo.
 - En el peor caso, la complejidad corresponde a la suma de las superficies de los bloques, aka $O(n)$.

1.3.3. Búsqueda por Interpolacion/Extrapolacion

1. Introduccion:

- Haria busqueda binaria en un guia telefonico para el nombre “Barbay”? En un diccionario para la ciudad “Zanzibar”?
- ojala que no: se puede provechar de la informacion que da la primera letra de cada palabra

2. Algoritmo

- Interaccion

3. Analisis * La analisis **en promedio** es complicada: conversamos solamente la intuicion matematica (para mas ver la publicacion cientifica de SODA04, Demaine Jones y Patrascu):

- si las llaves son **distribuidas uniformemente**, la distancia en promedio de la posición calculada por interpolación **lineal** hasta la posición real es de $\sqrt{r-l}$.
- entonces, se puede reducir el tamaño del subarreglo de n a \sqrt{n} cada (dos) comparaciones
- la búsqueda por interpolación
 - en promedio,
 - si las llaves son **distribuidas uniformemente**,
 - toma $O(\lg \lg n)$ comparaciones

4. Interacción.

5. Variantas

a) Interpolación non-lineal

- en un anuario telefónico o en un diccionario, las frecuencias de las letras **no** son uniformes

b) Búsqueda por Interpolación Mixta con Binaria

- Se puede buscar en tiempo
 - $O(\lg n)$ en el peor caso Y
 - $O(\lg \lg n)$ en el caso promedio?
- Solución fácil
- Solución más compleja

c) Búsqueda por Extrapolación

- Tarea 3

d) Búsqueda por Extrapolación Mixta con Doblada

- Tarea 3

6. Discussion:

- Porque todavía estudiar la complejidad en el modelo de comparaciones?
 - Cuando el peor caso es importante
 - Cuando la distribución no es uniforme o no es conocida
 - cuando el costo de la evaluación es más costoso que una simple comparación (en particular para la interpolación non lineal)

1.3.4. Estructuras de Árboles con Dominio Discreto

Tries o Árboles Digitales Ordenamos usualmente como pre-computación para buscar después. En el caso donde n es demasiado grande, ordenar puede ser demasiado caro. Consideramos alternativas para buscar.

1. Ejemplo de trie

- Insertar los nodos siguiente en un trie
 - hola
 - holístico
 - holograma
 - hologramas
 - ola
 - ole

2. Búsqueda

- con arreglos de tamaño σ en cada nodo:
 - $O(l)$ tiempo, pero $O(L\sigma)$ espacio
- con arreglos de tamaño variables en cada nodo:
 - $O(l \lg \sigma)$ tiempo (búsqueda binaria), $O(L)$ espacio (óptima).
- con hashing
 - $O(l)$ tiempo en promedio, $O(L)$ espacio.

3. Inserción

- Insertar “hora” en el árbol precedente
- Insertar “holístico” en el árbol precedente
- Borrar “hola” y “holística”
 - (TAREA)
 - Tiene que “limpiar”, pero no costo más que un factor constante de la búsqueda.

4. Espacio en el peor caso

- Pero caso por variable fijas:
 - n , la cantidad de llaves en el diccionario
 - σ , el alfabeto
 - m , el tamaño máximo de una llave
 - (σ^m llaves posibles).
- $\log_\sigma n$ niveles donde los arreglos son llenos
- $m - \log_s \text{igman}$ niveles donde cada arreglo contiene solamente un puntero.
- El costo total es del orden de $O((n + m)\sigma)$
 - $\sigma \times$ la cantidad de nodos.
 - los $\log_s \text{igman}$ primeros niveles se suman a $\$1 + \sigma + \sigma^2 + \dots \$$ nodos, que son $\frac{n\sigma - 1}{\sigma - 1}$
 - los otros niveles costarán $n(m - \log_s \text{igman})\sigma$
 - En total son $O(n\sigma(1 + m - \log_s \text{igman}))$

5. BONUS: PAT Trie

- Comprime las ramas de nodos de grado uno en una sola arista.
- La cadena (“string”) etiquetando la arista se guarda en el nodo hijo de la arista.
- Superior tanto en tiempo como en espacio en práctica.

Árboles de Sufijos (y Arreglos de Sufijos)

1. Árbol de Sufijos

- Espacio $O(n)$
- Construcción $O(n)$
- Búsqueda $O(m)$
- expresión regular $O(n^\lambda)$ donde $0 \leq \lambda \leq 1$

2. Arreglo de Sufijos

- Lista de sufijos ordenados
- cada busqueda de patrones costa dos busquedas binarias, donde cada comparacion costa $\leq m$, resultando en una complejidad de $O(m \lg n)$

3. BONUS: Rank en Bitmaps

- $\text{rank}(B, i)$ = cantidad de unos en $B[1, i]$
- consideramos
 - B estatico
 - se puede almacenar $\lg n$ bits.
- Solucion de Munro, Raman y Raman:
 - $b = 1/2 \lg n$ y $s = \lg^2 n$
 - Dividimos el index de B en
 - s Superbloques de tamaño $n/s \lg n$ bits
 - b Mini bloques de tamaño $n/b \lg s$ bits

$$\frac{n}{1/2 \lg n} \lg(\lg^2 n) = \frac{4n \lg \lg n}{\lg n} \in o(n)$$

- un diccionario con todos los bit vectores de tamaño

$$\sqrt{n} \lg n / 2 \lg \lg n \in o(n)$$

1.3.5. Hashing

Introduccion * Motivaciones

- Mejor tiempo **en promedio**
- Uso de todas la herramientas que tenemos
 - dominio de los valores
 - distribuciones de probabilidades de los valores

* Terminologia

- Tabla de Hash
 - Arreglo de tamaño N
 - que contiene n elementos a dentro de un universo $[1..U]$
- Funcion de Hash $h(K)$
 - $h : [1..U] \rightarrow [0..N - 1]$
 - se calcula rapidamente
 - distribue uniformemente (mas o menos) las llaves en la tabla en el caso ideal, $P[h(K) = i] = 1/N, \forall K, i$
- Collision
 - cuando $h(K_1) = h(K_2)$

Paradoxe del cumpleaños (Probabilidad de Colision)

- la probabilidad de collision es alta: “Paradoxe del cumpleaños”
 - (<http://www.wimp.com/justcoincidence/http://www.wimp.com/justcoincidence/>)
- Cual es la probabilidad que en una pieca de n personas, dos tiene la misma fecha de cumpleaños (a dentro de 365 dias)?
 - probabilidad que cada cumpleaños es unico:

$$\frac{364!}{(365 - n)! \times 365^{n-1}}$$

- Probabilidad que hay al menos un cumpleaños compartido:

$$\frac{1 - 364!}{(365 - n)! \times 365^{n-1}}$$

n	Proba
10	.12
23	.5
50	.97
100	.9999996

Hashing Abierto

1. Idea principal:

- resolver las colisions con caldenas

2. Ejemplo: (muy irealistico)

- $h(K) = K \bmod 10$
- Secuencia de insercion 52, 18, 70, 22, 44, 38, 62
 - Insertando al final (si hay que probar por repeticiones)

0	70
1	
2	52,22,62
3	
4	44
5	
6	
7	
8	18,38
9	

- Insertando al final (si no hay que probar por repeticiones)

0	70
1	
2	62,22,52
3	
4	44
5	
6	
7	
8	38,18
9	

3. Analisis:

- factor de carga es $\lambda = \frac{n}{N}$
- Rendimiento en el peor caso: $O(n)$

Hashing Cerrado

1. Idea principal:

- resolver las colisiones con busqueda, i.e.
 - $(h(K) + f(i)) \bmod N$
- diferentes tipos de busqueda:
 - **lineal** $f(i) = i$
 - primary “clustering” (formacion de secuencias largas)
 - **cuadratica** $f(i) = i^2$
 - secundario “clustering” (si $h(K_1) = h(K_2)$, la secuencias son las mismas.
 - **dobles hashing** $f(i) = i \cdot h'(K)$
 - $h'(K)$ debe ser prima con N

Ideal Hashing

- Imagina una funcion de hash que genera una secuencia que parece aleatoria.
- cada posicion tiene la misma probabilidad de ser la proxima
 - Probabilidad λ de elegir una posicion ocupada
 - Probabilidad $1 - \lambda$ de elegir una posicion libre
 - la secuencia de prueba puede tocar la misma posicion mas que una vez.
 - llaves identicas todavia siguen la misma secuencia.
- Cual es el costo promedio u_j de una busqueda negativa con j llaves en la tabla?
 - $\lambda = \frac{j}{N}$
 - $u_j = 1(1 - \lambda) + 2\lambda(1 - \lambda) + r\lambda^2(1 - \lambda) + \dots$
 - $= 1 + \lambda + \lambda^2 + \dots$
 - $= \frac{1}{1 - \lambda}$
 - $= \frac{1}{1 - j/N}$
 - $= \frac{N}{N - j} \in [1..N]$
- Cual es el costo promedio de una busqueda positiva s_i para el i -th elemento insertado?
 - $s_i = \frac{1}{1 - i/N} = \frac{N}{N - i}$
 - $s_n = 1/n \sum \frac{N}{N - i}$
 - $= \frac{N}{n} \sum_{i=0}^{n-1} \frac{1}{N - i}$
 - $= \frac{1}{\alpha} \sum_{i=0}^{n-1} \frac{1}{N - i}$ con $\alpha = \frac{n}{N}$
 - $< \frac{1}{\alpha} \int_{N-n}^N \frac{1}{x} dx$
 - $= 1/\alpha \ln \frac{N}{N-n}$
 - $= 1/\alpha \ln \frac{1}{\alpha}$
- Para $\alpha = 1/2$, el costo promedio es 1,387
- Para $\alpha = 0,9$, el costo promedio es 2,559

Universal Hashing

- $h(K) = ((aK + b) \bmod p) \bmod N$
- $a \in [1..p-1]$ elegido al azar
- $b \in [0..p-1]$ elegido al azar
- p es primo y mas grande que N
- N no es necesariamente primo

1.3.6. Hashing en memoria externa

1. Que pasa si la tabla de hashing no queda en memoria?

- IDEA: Simula un B-arbol de altura dos
 - organiza el dato con valores de hash
 - guarda un index en el nodo raiz
 - usa solamente **una parte** de la valor de hash para elegir el sobre-arbol
 - extiende el index cuando mas dato es agregado

2. Ideas de soluciones

- usar tecnicas similares a B -arboles, vEB arboles.

3. Descripcion

- B - cantidad de elementos en una pagina
- h - funcion de hash $\rightarrow [0..2^k - 1]$
- D - **profundidad general**, con $D \leq k$
 - la raiz tiene 2^D punteros a las paginas horas
 - la raiz es indexada con los D primeros bits de cada valor de hash.
- d_l - **profundidad local** de cada hora l
 - Las valores de hash en l tienen en comun los primeros d_l bits.
 - Hay 2^{D-d_l} punteros a la hora l
 - Siempre, $d_l \leq D$

4. Ejemplo

- $B = 4, k = 6, D = 2$

$d_l = 2$ 000100
 001000
 001011
 001100

$d_l = 2$ 010101
 011100

$d_l = 1$ 100100
 101101
 110001
 111100

5. Algoritmos

- Buscar
- Insertar
- Remover

6. Analisis

- Buscar, insertar remover
 - 1 acceso a la memoria secundariaa si el index se queda
- cantidad Promedio de paginas para tener n llaves
 - $\frac{n}{B \lg 2} \approx 1,44 \frac{n}{B}$
 - paginas son llenas a 69 % mas o menos.

1.4. Tecnicas de Analisis

1.4.1. Introduccion

* Cual es la performancia en el peor caso de las tecnicas en dominios discretos?

- Peor caso sobre
 - los datos, o
 - las consultas?
- Varios Problemas
 - ordenamiento
 - busqueda en arreglos ordenados
 - busqueda desordenada (hashing)

* Cual otra tecnica de analisis se puede usar para analizar la performancia de estas estructuras de datos y algoritmos?

- analisis amortizada
- analisis adaptativa
- analisis competitiva (de algoritmos en linea)

1.4.2. Analisi amortizada

MATERIAL A LEER

- “Amortized Analysis Explained” by Rebecca Fiebrink
- <http://www.cs.princeton.edu/fiebrink/423/AmortizedAnalysisExplainedFiebrink.pdf> [http : //www.cs.princeton.edu/fiebrink/423/AmortizedAnalysisExplainedFiebrink.pdf](http://www.cs.princeton.edu/fiebrink/423/AmortizedAnalysisExplainedFiebrink.pdf)
- Amortized Analysis
 - CLRS, Chapter 17: Amortized Analisis p.405-430
- Árbol biselado (Splay Trees)
 - <http://es.wikipedia.org/wiki/>

Principio de Analisis amortizada * Costo Amortizado:

- Se tiene una secuencia de n operaciones con costos c_1, c_2, \dots, c_n .
- Se quiere determinar $C = \sum_{i \in [1..n]} c_i$.
- Se puede tomar el peor caso de $c_i \leq t$ para tener una cota superior de $C \leq tn$.
- Un mejor analisis puede analizar el costo amortizado, con varias tecnicas:
 - analisis agregada
 - contabilidad de costos
 - funcion potencial.

* Aplicaciones:

- Move To Front
- “Self-adjusting and balanced binary trees” (Splay Trees”
- union-find data-structures
- max flow
- Fibonacci heaps
- dynamic array (vector en Java)

* Tres tecnicas basicas:

1. “*Aggregate analysis*”

- bound las proporciones de operaciones de cada tipo
- (e.g. mas inserciones que deletiones)

2. “*Accounting Method*”

- asigna un costo (positivo o negativo) a cada operacion
- ejemplos:
 - stack
 - java vector
- costo amortizado es la suma de los costos.

3. “*Potential Method*” (CLRS p.405)

- asigna una funcion de “energia potencial” (como en fisica)
 - (accounting method = height, potential method = potential energy)
- costo amortizado de una operacion es su costo mas el cambio de funcion de potencial que resulta de la operacion.
- es suficiente de asegurarse que la funcion de potencial es siempre mas grande que su valor inicial para mostrar que el costo total amortizado es una cota superior sobre el costo total de las operaciones.
- ejemplo:
 - analisis del algoritmo para min y max
 - analisis de MTF

* **Ejemplo: Incremento binario**

- Incrementar n veces un numero binario de k bits,
- e.g. desde cero hasta $2^k - 1$, con $n = 2^k$.
- costo $\leq kn$ (brute force)
- costo $\leq n + n/2 + \dots \leq 2n$ (costo amortizado)
- La tecnica usada aqui es la contabilidad de costos:
 - un flip de 0 a 1 cuesta 2
 - un flip de 1 a 0 cuesta 0
 - cada incremento cuesta 2.
- Analisis
 - ϕ = cantidad de unos en el numero
 - $\phi_0 = 0$
 - $c_i = l + 1$ cuando hay l unos
 - $\Delta\phi_i = -l + 1$
 - $\sum \bar{c}_i = \sum c_i + \phi_n - \phi_0 \geq 2$

* **Ejemplo: arreglo dinamico**, e.g. java Vector

- considera el tipo “Vector” en Java.
- de tamaño fijo n
- cuando accede a $n + 1$, crea un otro arreglo de tamaño $2n$, y copia todo.
- cual es el costo amortizado si agregando elementos uno a uno?

* **Ejemplo: Move-to-Front**

- analisis amortizada se puede usar para mostrar que MTF siempre performa a dentro de un factor de 4 de cualquier algoritmo (incluido un algoritmo optimal que conoce la secuencia de busquedas desde el inicio).
- Foncion de potencial es $2 \times \text{Inv}(\text{MTF})$

* **Ejemplo: Splay Arboles**

- el costo amortizado de cada insercion es $O(\lg n)$.

* **Ejemplo: Fibonacci Heap**

Analisi amortizada de Colas de Prioridades * Problema: Dado un conjunto (dinamica) de n tareas con valores, elegir y remudar la tarea de valor maxima.

* operaciones basicas:

- $M.\text{build}(\{e_1, e_2, \dots, e_n\})$
- $M.\text{insert}(e)$
- $M.\text{min}$
- $M.\text{deleteMin}$

* operaciones adicionales (“Addressable priority queues”)

- $M.insert(e)$, volviendo un puntero h (“handle”) al elemento insertado
- $M.remove(h)$, remudando el elemento especificado para h
- $M.decreaseKey(h,k)$, reduciendo la llave del elemento especificado para h
- $M.merge(Q)$, agregando el heap Q al heap M .

* Soluciones (conocidas o no):

	Linked List	Binary Tree	(Min-)Heap	Fibonacci Heap	Brodal Queue ¹
insert	$O(1)$	$O(\lg n)$	$O(\lg n)$	$O(1)$	$O(1)$
accessmin	$O(n)$	$O(1)$	$O(1)$	$O(1)$	$O(1)$
deletemin	$O(n)$	$O(\lg n)$	$O(\lg n)$	$O(\lg n)^*$	$O(\lg n)$
decreasekey	$O(1)$	$O(\lg n)$	$O(\lg n)$	$O(1)^*$	$O(1)$
delete	$O(n)$	$O(n)$	$O(\lg n)$	$O(\lg n)^*$	$O(\lg n)$
merge	$O(1)$	$O(m \lg(n + m))$	$O(m \lg(n + m))$	$O(1)$	$O(1)$

1. Colas de prioridades binarias (“Binary Heaps”)

* La solución tradicional

- un arreglo de n elementos
- hijos del nodo i en posiciones $2i$ y $2i + 1$
- la valor de cada nodo es mas pequena que las valores de su hijos.

2. Cuidado de no implementar $M.build(\{e_1, e_2, \dots, e_n\})$ con n inserciones (sift up $\rightarrow O(n \lg n)$), pero con $n/2$ sift-down ($\rightarrow O(n)$).

3. En $M.deleteMin()$, algunas variantes de implementación (después de cambiar el min con $A[n]$):

- a) dos comparaciones en cada nivel hasta encontrar la posición final de $A[n]$
 - $2 \lg n$ comparaciones en el peor caso
 - $\lg n$ copias
- b) una comparación en cada nivel para definir un camino de tamaño $\lg n$, y una búsqueda binaria para encontrar la posición final
 - $\lg n + O(\lg \lg n)$ comparaciones en el peor caso
 - $\lg n$ copias en el peor caso
- c) una comparación en cada nivel para definir un camino de tamaño $\lg n$, y una búsqueda **secuencial** up para encontrar la posición final
 - $2 \lg n$ comparaciones en el peor caso
 - $\lg n$ copias en el peor caso
 - pero en práctica y promedio mucho mejor.

4. Porque la complejidad es mejor con un heap que con un arreglo ordenado?

- Para cada conjunto, hay solamente un arreglo ordenado que puede representarlo, pero muchos heaps posibles: eso da mas flexibilidad para la mantención dinámica de la estructura de datos.

5. Colas de prioridades con punteros (“Addressable priority queues”)

- Algun que mas flexible que los arreglos ordenados, las colas de prioridades binarias todavia son de estructura muy estricta, por ejemplo para la union de filas. Estructuras mas flexibles consideran un “bosque” de arboles. Ademas, estas estructuras de arboles son implementadas con punteros (en ves de implementarlos en arreglos).
 - Hay diferentes variantes. Todas tienen en comun los puntos siguientes:
 - a) un puntero **minPtr** indica el nodo de valor minima en el bosque, raiz de alguno arbol.
 - b) **insert** agregas un nuevo arbol al bosque en tiempo $O(1)$
 - c) **deleteMin** remudas el nodo indicado par minPtr, dividiendo su arbol en dos nuevos arboles. Buscamos para el nuevo min y fusionamos algunos arboles (los detalles diferencian las variantes)
 - d) **decreaseKey(h,k)** es implementado cortando el arbol al nodo indicado por h, y rebalanceando el arbol cortado.
 - e) **delete()** es reducido a **decreaseKey(h,0)** y **deleteMin**
 - f) “Pairing Heaps”
 - Malo rendimiento en el peor caso, pero bastante buena en
6. rebalancea los arboles solamente en deleteMin, y solamente con pares de raises (i.e. la cantidad de arboles es reducida por dos a cada deleteMin).

Operacion	Amortizado
Insert(C,x)	$O(1)$
Merge	$O(1)$
ExtractMin	$O(\lg n)$
decreaseKey(h,k)	$\Omega(n \lg n \lg \lg n)$

7. Colas de prioridades binomiales (“Binomial Heaps”)

http://en.wikipedia.org/wiki/Binomial_heap[http : // en . wikipedia . org / wiki / Binomial_heap](http://en.wikipedia.org/wiki/Binomial_heap)

* Definicion

- Un **arbol binomial** (de orden k) tiene exactamente k hijos de orden distintos $k - 1, k - 2, \dots, 0$. [Un arbol binomial de orden 0 tiene 0 hijos.]
- Un **bosque binomial** es un conjunto de arboles binomiales de orden **distintas** (i.e. hay cero o uno arboles de cada orden).

8. Para cada arbol T

- $h(T) \leq \lg|T|$
- $|T| \geq 2^{h(T)}$

9. Para el bosque

- $\forall n$ hay solamente uno bosque binomial con n nodos.
- al maxima tiene $\lfloor \lg(n+1) \rfloor$ arboles.
- la descomposicion del bosque en arboles de orden k corresponde a la descomposicion de n en base de dos.

10. una **cola binomial** es un bosque binomial donde cada nodo almacena una clave, y siempre la clave de un padre es inferior o igual a la clave de un hijo.

* Operaciones

Operacion	Peor Caso
Merge	$O(\lg n)$
FindMin	$O(\lg n)$
ExtractMin	$O(\lg n)$
Insert(C,x)	$O(\lg n)$
Heapify	$O(n)$
remove(h)	$O(\lg n)$
decreaseKey(h,k)	$O(\lg n)$
merge(Q)	$O(\lg n)$

* Union

11. Union de dos arboles binomiales de mismo orden:

- agrega T_2 a T_1 si T_1 tiene la raiz mas pequena.

12. Union de dos bosques binomiales:

- si hay uno arbol de orden k , es lo de la union
- si hay dos arboles de orden k , calcula la union en un arbol de orden $k + 1$
- la propagacion es similar a la suma de enteros en binario.

13. Complejidad

- $O(\lg n)$ en el peor caso

14. agrega un arbol de orden 0 y hace la union si necesitado

15. Complejidad $O(\lg n)$ en el peor caso

16. Puede ser $O(1)$ sin corregir el bosque, que tiene de ser corregido mas tarde, que puede ser en tiempo $O(n)$ peor caso, pero sera $O(\lg n)$ en tiempo amortizado.

* Minima

17. lei la lista de al maxima $\lfloor \lg(n + 1) \rfloor$ raices

18. Complejidad $O(\lg n)$

19. Puede ser $O(1)$ si precalculando un puntero al minima, que tiene de ser corregido (en tiempo $O(\lg n)$) a cada modificacion.

* DeleteMin

20. encontra el min

21. remuda el min de su arbol (la raiz)

22. reordena su hijos para su orden, en un bosque binomial

23. hace la union con el bosque binomial original, menos el arbol del min

24. complejidad $O(\lg n)$

* DecreaseKey

25. sigue el camino abajo hasta que la condicion del heap es corregida.

26. cada arbol tiene altura $\lg n$, entonces la complejidad es $O(\lg n)$ en el peor caso.

* Delete

Finger Search Tree: la busqueda doblada de los arboles de busqueda

APUNTES

- Estructura de datos
- algoritmo de búsqueda
- análisis: búsqueda en $O(\lg p)$

Algoritmo de Ordenamiento adaptivos basicos

Merge Sort Adaptivo: Runs

Local Insertion Sort: Inv

Another Insertion Sort: REM

Computacion de la Union

Computacion de la Interseccion

1.4.4. Algoritmos en linea

List Accessing REFERENCIA: Capitulo 1 en “Online Computation and Competitive Analysis”, de Allan Borodin y Ran El-Yaniv

- “List Accessing”
 - Considera la secuencia de búsqueda de tamaño n , en un diccionario de tamaño σ : “1,1,1,1,1,2,2,2,3,3,3,4,4,4,5,5,5,6,6,6,7,7,7,8,8,8,9,9,9,10,10,10,11,11,11,12,12,12,13,13,13,14,14,14,15,15,15,16,16,16,17,17,17,18,18,18,19,19,19,20,20,20,21,21,21,22,22,22,23,23,23,24,24,24,25,25,25,26,26,26,27,27,27,28,28,28,29,29,29,30,30,30,31,31,31,32,32,32,33,33,33,34,34,34,35,35,35,36,36,36,37,37,37,38,38,38,39,39,39,40,40,40,41,41,41,42,42,42,43,43,43,44,44,44,45,45,45,46,46,46,47,47,47,48,48,48,49,49,49,50,50,50,51,51,51,52,52,52,53,53,53,54,54,54,55,55,55,56,56,56,57,57,57,58,58,58,59,59,59,60,60,60,61,61,61,62,62,62,63,63,63,64,64,64,65,65,65,66,66,66,67,67,67,68,68,68,69,69,69,70,70,70,71,71,71,72,72,72,73,73,73,74,74,74,75,75,75,76,76,76,77,77,77,78,78,78,79,79,79,80,80,80,81,81,81,82,82,82,83,83,83,84,84,84,85,85,85,86,86,86,87,87,87,88,88,88,89,89,89,90,90,90,91,91,91,92,92,92,93,93,93,94,94,94,95,95,95,96,96,96,97,97,97,98,98,98,99,99,99,100,100,100,101,101,101,102,102,102,103,103,103,104,104,104,105,105,105,106,106,106,107,107,107,108,108,108,109,109,109,110,110,110,111,111,111,112,112,112,113,113,113,114,114,114,115,115,115,116,116,116,117,117,117,118,118,118,119,119,119,120,120,120,121,121,121,122,122,122,123,123,123,124,124,124,125,125,125,126,126,126,127,127,127,128,128,128,129,129,129,130,130,130,131,131,131,132,132,132,133,133,133,134,134,134,135,135,135,136,136,136,137,137,137,138,138,138,139,139,139,140,140,140,141,141,141,142,142,142,143,143,143,144,144,144,145,145,145,146,146,146,147,147,147,148,148,148,149,149,149,150,150,150,151,151,151,152,152,152,153,153,153,154,154,154,155,155,155,156,156,156,157,157,157,158,158,158,159,159,159,160,160,160,161,161,161,162,162,162,163,163,163,164,164,164,165,165,165,166,166,166,167,167,167,168,168,168,169,169,169,170,170,170,171,171,171,172,172,172,173,173,173,174,174,174,175,175,175,176,176,176,177,177,177,178,178,178,179,179,179,180,180,180,181,181,181,182,182,182,183,183,183,184,184,184,185,185,185,186,186,186,187,187,187,188,188,188,189,189,189,190,190,190,191,191,191,192,192,192,193,193,193,194,194,194,195,195,195,196,196,196,197,197,197,198,198,198,199,199,199,200,200,200,201,201,201,202,202,202,203,203,203,204,204,204,205,205,205,206,206,206,207,207,207,208,208,208,209,209,209,210,210,210,211,211,211,212,212,212,213,213,213,214,214,214,215,215,215,216,216,216,217,217,217,218,218,218,219,219,219,220,220,220,221,221,221,222,222,222,223,223,223,224,224,224,225,225,225,226,226,226,227,227,227,228,228,228,229,229,229,230,230,230,231,231,231,232,232,232,233,233,233,234,234,234,235,235,235,236,236,236,237,237,237,238,238,238,239,239,239,240,240,240,241,241,241,242,242,242,243,243,243,244,244,244,245,245,245,246,246,246,247,247,247,248,248,248,249,249,249,250,250,250,251,251,251,252,252,252,253,253,253,254,254,254,255,255,255,256,256,256,257,257,257,258,258,258,259,259,259,260,260,260,261,261,261,262,262,262,263,263,263,264,264,264,265,265,265,266,266,266,267,267,267,268,268,268,269,269,269,270,270,270,271,271,271,272,272,272,273,273,273,274,274,274,275,275,275,276,276,276,277,277,277,278,278,278,279,279,279,280,280,280,281,281,281,282,282,282,283,283,283,284,284,284,285,285,285,286,286,286,287,287,287,288,288,288,289,289,289,290,290,290,291,291,291,292,292,292,293,293,293,294,294,294,295,295,295,296,296,296,297,297,297,298,298,298,299,299,299,300,300,300,301,301,301,302,302,302,303,303,303,304,304,304,305,305,305,306,306,306,307,307,307,308,308,308,309,309,309,310,310,310,311,311,311,312,312,312,313,313,313,314,314,314,315,315,315,316,316,316,317,317,317,318,318,318,319,319,319,320,320,320,321,321,321,322,322,322,323,323,323,324,324,324,325,325,325,326,326,326,327,327,327,328,328,328,329,329,329,330,330,330,331,331,331,332,332,332,333,333,333,334,334,334,335,335,335,336,336,336,337,337,337,338,338,338,339,339,339,340,340,340,341,341,341,342,342,342,343,343,343,344,344,344,345,345,345,346,346,346,347,347,347,348,348,348,349,349,349,350,350,350,351,351,351,352,352,352,353,353,353,354,354,354,355,355,355,356,356,356,357,357,357,358,358,358,359,359,359,360,360,360,361,361,361,362,362,362,363,363,363,3

- Estos son “Algoritmos en Linea”

- algoritmo de optimizacion
- que conoce solamente una parte de la entrada al tiempo t .
- se compara a la competitividad con el algoritmo offline que conoce toda la instancia.
- Como se puede medir su complejidad?
 - cada algoritmo ejecuta $O(n)$ comparaciones para cada busqueda en el peor caso!!!!
 - tiene de considerar instancias “faciles” y “dificiles”
 - una medida de dificultad
 - e.g. el rendimiento del *mejor algoritmo “offline”*

- Competitive Analysis: instancias “dificiles” o “faciles”

- Las estructuras de datos dinamicas pueden aprovechar de secuencias “faciles” de consultas: eso se llama “online”.
- pero para muchos problemas online, todas las heurísticas se comportan de la misma manera en el peor caso.
- Por eso se identifica una medida de dificultad de las instancias, y se comparan los rendimientos de los algoritmos sobre instancias que tienen una valor fijada de este medida de dificultad.
- Tradicionalmente, esta medida de dificultad es el rendimiento del mejor algoritmo “offline”: eso se llama **competitive analysis**, resultando en el **competitive ratio**, el ratio entre la complejidad del algoritmo ONLINE y la complejidad del mejor algoritmo OFFLINE.
 - por ejemplo, veamos que MTF tiene un competitive ratio de 2
- Pero todavia hay algoritmos con performance practicas muy distintas que tienen el mismo competitive ratio. Por eso se introduce otras medidas de dificultades mas sofisticadas, y mas especialidades en cada problema.

- Competitividad

* Optimizacion/aproximacion

- A es $k(n)$ **competitiva** para un problema de **minimizacion** si

$$\exists b \forall n, x, |x| = n, C_A(x) - k(n)C_{OPT}(x) \leq b$$

- A es $k(n)$ **competitiva** para un problema de **maximizacion** si

$$\exists b, \forall n, x, |x| = n, C_{OPT}(x) - k(n)C_A(x) \leq b$$

- an algoritmo en linea es **c-competitiva** si

$$\exists \alpha, \forall I ALG(I) \leq cOPT(I) + \alpha$$

- an algoritmo en linea es **estrictamente c-competitiva** si

$$\forall I ALG(I) \leq cOPT(I)$$

- Sleator-Tarjan sobre MTF

- costo de una busqueda negativa (la llave NO esta en el diccionario)
 - σ
- costo de una busqueda positiva (la llave esta en el diccionario)

- la posición de la llave, no más que σ
- en promedio para una distribución de probabilidad fijada:
 - ◊ $MTF \leq 2OPT$,

◦ Prueba: (from my notes in my CS240 slides)

How does MTF compare to the optimal ordering?

- ◊ Assume that:
 - ◊ the keys k_1, \dots, k_n have probabilities $p_1 \geq p_2 \geq \dots \geq p_n \geq 0$
 - ◊ the list is used sufficiently to reach a steady state.
- ◊ Then:

$$C_{MTF} < 2 \cdot C_{OPT}$$

- ◊ Proof:
 - ◊ $C_{OPT} = \sum_{j=1}^n j p_j$
 - ◊ $C_{MTF} = \sum_{j=1}^n p_j (\text{cost of finding } k_j)$
 - ◊ $C_{MTF} = \sum_{j=1}^n p_j (1 + \text{number of keys before } k_j)$
 - ◊ To compute the average number of keys before k_j :

$$\Pr[k_i \text{ before } k_j] = \frac{p_i}{p_i + p_j}$$

$$E(\text{number of keys before } k_j) = \sum_{i \neq j} \frac{p_i}{p_i + p_j}$$

- ◊ k_i is before k_j if and only if k_i was accessed more recently than k_j .
- ◊ Consider the last time either k_i or k_j was looked up. What is the probability that it was k_i ?

$$P(k_i \text{ before } k_j) = P(k_i \text{ chosen} \mid k_i \text{ or } k_j \text{ chosen})$$

$$P(k_i \text{ before } k_j) = \frac{P(k_i \text{ chosen})}{P(k_i \text{ or } k_j \text{ chosen})}$$

$$P(k_i \text{ before } k_j) = \frac{p_i}{p_i + p_j}$$

- ◊ Therefore,
- ◊ Joining both previous formulas:

$$C_{MTF} = \sum_{j=1}^n p_j \left(1 + \sum_{i \neq j} \frac{p_i}{p_i + p_j} \right)$$

- ◊ reordering the terms:

$$C_{MTF} = 1 + 2 \sum_{j=1}^n \sum_{i < j} \frac{p_i p_j}{p_i + p_j}$$

- ◊ Because $\frac{p_i}{p_i + p_j} \leq 1$:

$$C_{MTF} \leq 1 + 2 \sum_{j=1}^n p_j \left(\sum_{i < j} 1 \right)$$

$$C_{MTF} = 1 + 2 \sum_{j=1}^n p_j (j - 1)$$

$$C_{MTF} = 1 + 2C_{OPT} + 2 \sum_{j=1}^n (-p_j)$$

◇ Because $\sum_{j=1}^n (p_j) = 1$:

$$C_{MTF} = 2C_{OPT} - 1$$

BONUS Aplicaciones a la compression de textos

* Bentley, Sleator, Tarjan and Wei proponieron de comprimir un texto utilizando una lista dinamica, donde el codigo para un simbolo es la posicion del simbolo en la lista.

- Experimentalmente, se compara a Huffman:

- a veces mucho mejor
- nunca mucho peor.

- Experimentalmente, 6% mejor que GZip, que es enorme!

Paginamiento Deterministico REFERENCIA: Capitulo 2 en “Online Computation and Competitive Analysis”, de Allan Borodin y Ran El-Yaniv

1. Paginamiento

- Definicion:

- elegir cual paginas guardar en memoria, dado
 - una secuencia online de n consultas para paginas, y
 - un cache de k paginas.

- Politicas:

- LRU (Least Recently Used)
- CLOCK (1bit LRU)
- FIFO (First In First Out)
- LFU (Least Frequently Used)
- LIFO = MRU (Most Recently Used)
- FWF (Flush When Full)
- LFD (Offline, Longest Forward Distance)

- Ustedes tienen una idea de cuales son las peores/mejores?

2. Relacion con “List Accessing”

a) Cada “List accessing” algoritmo corresponde a un algoritmo de paginamiento:

- cada miss, borra el ultimo elemento de la lista y “inserta” el nuevo elemento.

b) No hay una reduccion tan clara en la otra direccion.

3. Offline analysis

- LFD performa $O(n/k)$ misses
- Cualquier algoritmo Offline performa $\Omega(n/k)$ en el peor caso.

4. Online analysis: resultados basicos

a) $\forall A$ online, hay una entrada con n fallas.

- Estrategia de adversario.

b) No algoritmo online puede ser mejor que k competitivo.

- Obvio, comparando con LFD.

c) MRU=LIFO NO es competitivo

- Considera $S = p_1, p_2, \dots, p_k, p_{k+1}, p_k, p_{k+1}, p_k, p_{k+1}, p_k, \dots$
- despues las k primeras consultas, MRU va a tener un miss cada consulta, cuando LFD nunca mas.

d) LFU no es competitivo

- Considera $l > 0$ y $S = p_1^l, p_2^l, \dots, p_{k-1}^l, (p_k, p_{k+1})^{l-1}$
- Despues de las $(k-1)l$ primeras consultas, LFU va a tener un miss cada consulta, cuando LFD solamente dos.

e) Que tal de FWF?

- MRU=LIFO es un poco estúpido, su mala rendimiento no es una sorpresa.
- FWF es un algoritmo muy ingenuo tambien, pero vamos a ver que no tiene un rendimiento tan mal “en teoria”.

5. BONUS: Competitive Analysis: Algoritmos a Marcas

a) k -fases particiones

Para cada secuencia S , partitionala en secuencias S_1, \dots, S_δ tal que

- $S_0 = \emptyset$
- S_i es la secuencia Maxima despues de S_{i-1} que contiene al maximum k consultas distintas.
- Llamamos “fase i .” el tiempo que el algoritmo considera elementos de la subsecuencia S_i .
- Nota que eso es independiente del algoritmo considerado.

b) Algoritmo con marcas

- agrega a cada pagina de memoria lenta un bit de marca.
- al inicio de cada fase, remuda las marcas de cada pagina en memoria.
- a dentro de una fase, marca una pagina la primera vez que es consultada.

- un algoritmo a marca (“marking algorithm”) es un algoritmo que nunca remuda una pagina marcada de su cache.
- c) Un algoritmo con marcas es k -competitiva
- En cada fase,
 - un algoritmo ONLINE con marcas performa al maximum k miss.
 - Un algoritmo OFFLINE (e.g. LFD) performa al minimum 1 miss.
 - QED
- d) LRU, CLOCK y FWF son algoritmos con marcas
- e) LRU, CLOCK y FWF tienen un ratio competitivo OPTIMO

6. Mas resultados:

- a) La analisis se puede generalizar al caso donde el algoritmo offline tiene h paginas, y el algoritmo online tiene $k \geq h$ paginas.
- Cada algoritmo **con marcas** es $\frac{k}{k-h+1}$ -competitiva.
- b) Definicion de algoritmos (conservadores) da resultado similar para FIFO (que no es con marcas pero es conservador).
- c) En practica, sabemos que LRU es mucho mejor que FWF (for instancia). Habia mucha investigacion para intentar de mejorar la analisis por 20 años, ahora parece que hay una analisis que explica la mejor rendimiento de LRU sobre FWF, y de variantes de LRU que pueden saber x pasos en el futuro (Reza Dorrigiv y Alex Lopez-Ortiz).

BONUS: ”Ski Renting

http://en.wikipedia.org/wiki/Ski_rental_problem http://en.wikipedia.org/wiki/Ski_rental_problem

1.4.5. MAYB Complejidad Parametrizada

1.4.6. PREGUNTAS [0/1] :PREGUNTAS:

Analisis Amortizada: arreglo dinamico (Part 1) Queremos implementar una pila (“stack”) en un arreglo. Iniciamos con un arreglo de tamaño $s = 1$, y cuando se llena, creamos un arreglo mas grande, copiamos todo en en nuevo arreglo y sigamos.

Cual es el costo amortizado de una insercion si el nuevo arreglo es de tamaño $n + 1$?

1. $O(1)$
2. $O(\lg n)$
3. $O(n)$
4. $O(n^2)$
5. otra respuesta

Análisis Amortizada: arreglo dinámico (Part 2) ¿Cuál es el costo amortizado de una inserción si el nuevo arreglo es de tamaño $2n$?

1. $O(1)$
2. $O(\lg n)$
3. $O(n)$
4. $O(n^2)$
5. otra respuesta

TODO Análisis Amortizada: arreglo dinámico (Part 3) ¿Cuál es el costo amortizado de una inserción si el nuevo arreglo es de tamaño $4n$?

1. menos que 2
2. 2
3. entre 2 y 3
4. 3
5. más que 3

Análisis Amortizada: arreglo dinámico (Part 4) ¿Cuál es el costo amortizado de una inserción si el nuevo arreglo es de tamaño n^2 (y el primero arreglo de tamaño 2)?

1. $O(1)$
2. $O(\lg n)$
3. $O(n)$
4. $O(n^2)$
5. otra respuesta

1.5. Conclusion

1.6. RESUMEN Unidad 3

* Resultados de Aprendizajes de la Unidad

- Comprender las técnicas de algoritmos de
 - costo amortizado,
 - uso de finitud, y
 - algoritmos competitivos

- Ser capaz de diseñar y analizar algoritmos y estructuras de datos basados en estos principios.
- conocer algunos casos de estudio relevantes

* Principales casos de estudio:

- estructuras para union-find,

- colas binomiales
- splay trees,
- búsqueda por interpolación
- radix sort
- árboles de van Emde Boas
- árboles de sufijos
- técnica de los cuatro rusos,
- paginamiento
- búsqueda no acotada (unbounded search, doubling search)