

# Funciones de hash y ordenación lineal

Clase 13

IIC 2133 - Sección 2

Prof. Mario Droguett

# Sumario

**Introducción**

Funciones de hash

Ordenación lineal

Cierre

# Diccionarios

## Definición

Un **diccionario** es una estructura de datos con las siguientes operaciones

- **Asociar** un valor a una llave
- **Actualizar** el valor asociado a una llave
- **Obtener** el valor asociado a una llave
- En ciertos casos, **eliminar** de la estructura una asociación llave-valor

Los ABB fueron nuestra primera EDD para implementar diccionarios

# Funciones de hash

## Definición

Dado un espacio de llaves  $K$  y un natural  $m > 0$ , una **función de hash** se define como

$$h : K \rightarrow \{0, \dots, m-1\}$$

Dado  $k \in K$ , llamaremos **valor de hash de  $k$**  a la evaluación  $h(k)$ .

Notemos que

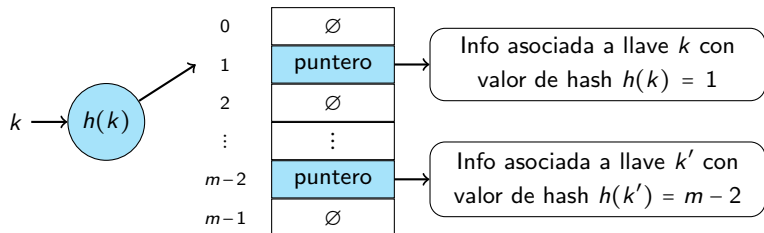
- Una función de hash nos permite mapear un espacio de llaves a otro más pequeño (con  $m$  razonable)
- Una función de hash no necesariamente es **inyectiva**
- Si  $m < |K|$ , entonces no puede ser inyectiva
- En la práctica,  $m \ll |K|$

# Tablas de hash

## Definición

Dado  $m > 0$  y un conjunto de llaves  $K$ , una **tabla de hash**  $T$  es una EDD que asocia valores a llaves indexadas usando una función de hash

$h: K \rightarrow \{0, \dots, m-1\}$ . Diremos que tal  $T$  es de tamaño  $m$ .



Una **colisión** ocurre cuando  $k_1 \neq k_2$  y  $h(k_1) = h(k_2)$

¿Cómo manejamos colisiones?

# Colisiones

Dos estrategias de manejo de colisiones

- Encadenamiento: listas ligadas en cada  $T[h(k)]$
- Direccionamiento abierto: mediante un tipo de sondeo, se insertan los datos en posiciones vacías en  $T$

Además, definimos una medida de qué tan llena está la tabla

## Definición

Dada una tabla de hash  $T$  de tamaño  $m$  con  $n$  valores almacenados, se define su **factor de carga** como

$$\lambda = \frac{n}{m}$$

Si  $\lambda$  no es aceptable, hacemos **rehashing**

# Objetivos de la clase

- ☐ Identificar buenas propiedades de una función de hash
- ☐ Comprender algunas familias de funciones de hash comunes
- ☐ Conocer algoritmos de ordenación en tiempo lineal
- ☐ Comprender la limitación de dominio para tener tales algoritmos

# Sumario

Introducción

**Funciones de hash**

Ordenación lineal

Cierre



# Funciones de hash

## Definición

Dado un espacio de llaves  $K$  y un natural  $m > 0$ , una **función de hash** se define como

$$h : K \rightarrow \{0, \dots, m-1\}$$

Dado  $k \in K$ , llamaremos **valor de hash de  $k$**  a la evaluación  $h(k)$ .

¿Qué caracteriza al espacio de llaves  $K$ ?

- En principio podemos aceptar cualquier tipo de conjunto  $K$
- Haremos la distinción entre llaves **numéricas** y **no numéricas**
- En cualquier caso, nos interesa llevar esas llaves a  $\{0, \dots, m-1\}$
- En general, las no numéricas primero se convierten a entero

Nos centraremos en las llaves numéricas

# Funciones de hash

Nos interesa tener *buenas* funciones de hash

- Esperamos que el mapeo de llaves sea uniforme en el rango  $\{0, \dots, m-1\}$
- Debe ser fácil/barata de calcular

¿Cómo lograr esto? ¿o al menos acercarnos?

Sabemos que las llaves de  $K$ , en última instancia, se representan en binario

- Hay operaciones en binario muy eficientes
- Idea: que el valor  $h(k)$  dependa de todos/muchos bits de  $k$

# Funciones de hash

## Ejemplo

Consideremos  $m = 16$  (tamaño de la tabla) y  $k = 1234567$ . El cuadrado de la llave es  $k^2 = 1524155677489$ , cuya representación binaria tiene  $n$  bits y es

$$(k^2)_2 = (010110001011011110110000011111101100110001)_2$$

Nos centraremos en los  $m = 16$  bits intermedios

$$01011000101101\underline{1111011000001111}1101100110001$$

Consideremos el valor de hash  $h(k)$  dado por

$$\underline{1111011000001111} = (62991)_2$$

de manera que  $h(1234567) = 62991$

# Funciones de hash

## Ejemplo

Generalicemos este resultado para definir una función de hash que considere los  $m$  bits centrales de  $k^2$

Recordemos que en representación binaria

- Dividir  $a$  por  $2^b$  **borra** los  $b$  bits de la derecha de  $(a)_2$
- $(2^b)_2$  es un 1 seguido de  $b$  ceros
- Como  $[a \bmod 2^b]$  es el resto al dividir entre  $2^b$  y  $[a \bmod 2^b] < 2^b$ , entonces  $(a \bmod 2^b)_2$  son los  $b$  bits de la derecha en  $(a/2^b)_2$

# Funciones de hash

## Ejemplo

La función de hash que definimos es

$$h_m(k) = \left( \frac{k^2}{2^r} \right) \bmod 2^m$$

donde  $r = n - m/2$  es la cant. de bits ignorados en el lado derecho de  $(k^2)_2$

- Modificando  $m$  podemos alterar la cantidad de bits que *importan*
- Si  $m = n$ , entonces  $r = 0$  y

$$h_n(k) = \left( \frac{k^2}{2^0} \right) \bmod 2^n = k^2 \bmod 2^n = k^2$$

que corresponde a la función  $h_m$  que más bits de  $k^2$  incorpora

# Hashing modular

Una familia simple de funciones razonable, que ya estudiamos, es de la forma

$$h(k) = k \bmod m$$

y se conoce como **hashing modular** o **método de la división**

- Es fácil de calcular
- Distribuye uniformemente las llaves si estas son *aleatorias*

# Hashing modular

Como se vio en el ejemplo, no siempre calculamos el módulo de un *k puro*... puede ser resultado de operaciones entre otros números

Recordemos las propiedades de la aritmética modular

- $(a + b) \bmod c = ((a \bmod c) + (b \bmod c)) \bmod c$
- $(a \cdot b) \bmod c = ((a \bmod c) \cdot (b \bmod c)) \bmod c$

## Ejemplo

Calculemos  $29 \cdot 72 \bmod 13$  sin efectuar el producto  $29 \cdot 72$

$$\begin{aligned} 29 \cdot 72 \bmod 13 &= ((29 \bmod 13) \cdot (72 \bmod 13)) \bmod 13 \\ &= (3 \cdot 7) \bmod 13 \\ &= 21 \bmod 13 \\ &= 8 \end{aligned}$$

# Hashing modular

## Algunas propiedades del hashing modular

- Si  $m$  es muy pequeño, habrán muchas colisiones
- Si  $m = 2^b$ , solo serán relevantes los  $b$  bits menos significativos... pueden repetirse mucho ciertos patrones
- Además, no queremos que  $m$  tenga muchos divisores
  - Si  $m$  es par y las llaves  $k$  son pares,  $k \bmod m$  es par
  - ¡La mitad de la tabla va a estar vacía!

Es conveniente tomar  $m$  como primo, cercano al tamaño ideal de la tabla



# Método de la multiplicación

Otra familia razonable considera un real  $0 < A < 1$  y es de la forma

$$h(k) = \lfloor m \cdot (A \cdot k \bmod 1) \rfloor$$

y se conoce como **método de la multiplicación**

- Se extrae la parte decimal de  $A \cdot k$
- A diferencia del hashing modular, el valor de  $m$  no es crítico
- Tomar  $m$  como potencia de 2 simplifica los cálculos

Es conveniente tomar  $m$  como primo, cercano al tamaño ideal de la tabla

# Llaves strings

Si las llaves son strings debemos convertirlas a enteros

Una forma de realizar la conversión a entero es usar códigos ASCII y concatenarlos

## Ejemplo

Si sabemos que

$$\text{ASCII}("A") = 64 \rightarrow 01000001 \quad \text{ASCII}("B") = 66 \rightarrow 01000010$$

entonces una posible interpretación numérica de  $s = "AB"$  es

$$0100000101000010 = (16706)_2$$

Ahora bien, si

$$s' = "A \text{ LONG EXAMPLE STRING}"$$

¿Qué tan grande sería esta representación para  $s'$ ?

# Llaves strings

Una forma común es interpretar  $s$  como entero en una **base apropiada**

Si  $s = s_0s_1 \dots s_p$ , una forma de convertirlo a un entero  $k$  es

$$k = \#(s_p) + \#(s_{p-1}) \cdot R + \dots + \#(s_0) \cdot R^p = \sum_{i=0}^p \#(s_{p-i}) \cdot R^i$$

donde

- $\#(s_i) = \text{ASCII}(s_i)$
- $R \in \{31, 37\}$ , i.e. un número primo que permite que los bits de todos los  $s_i$  jueguen un rol

Con esta estrategia estamos interpretando a  $s$  como un número de  $p + 1$  dígitos en base  $R$

Finalmente, convertimos  $k$  al rango  $\{0, \dots, m - 1\}$  con alguna función de hash

# Restricciones de dominio

El recorrido de una función de hash es de la forma  $\{0, \dots, m-1\}$

- Elementos del recorrido son naturales
- Esto permite usarlos para indexar posiciones de un arreglo
- Arreglo es de tamaño  $m$

¿Qué sucede si queremos ordenar un conjunto de naturales  $\{0, \dots, k\}$ ?

Podemos aprovechar que sirven como índices de un arreglo

# Sumario

Introducción

Funciones de hash

**Ordenación lineal**

Cierre

# Ordenación en tiempo lineal

Consideremos una  $A$  secuencia de  $n$  naturales entre 0 y  $k$

- Para todo  $a_i \in A$ , se tiene que  $0 \leq a_i \leq k$
- Notemos que no necesariamente  $n = k - 1$
- La secuencia  $A$  puede tener elementos repetidos

Propondremos un algoritmo de ordenación que **no compara**

- Para cada dato, contaremos cuántos datos son menores que él
- Esto nos indica la posición final de cada elemento

Si  $k \in \mathcal{O}(n)$ , entonces este algoritmo será  $\Theta(n)$

## El algoritmo CountingSort()

**input** : Arreglo  $A[0 \dots n-1]$ , natural  $k$

**output**: Arreglo  $B[0 \dots n-1]$

CountingSort ( $A, k$ ):

```
1    $B[0 \dots n-1] \leftarrow$  arreglo vacío de  $n$  celdas
2    $C[0 \dots k] \leftarrow$  arreglo vacío de  $k+1$  celdas
3   for  $i = 0 \dots k$  :
4        $C[i] \leftarrow 0$ 
5   for  $j = 0 \dots n-1$  :
6        $C[A[j]] \leftarrow C[A[j]] + 1$ 
7   for  $p = 1 \dots k$  :
8        $C[p] \leftarrow C[p] + C[p-1]$ 
9   for  $r = n-1 \dots 0$  :
10       $B[C[A[r]] - 1] \leftarrow A[r]$ 
11       $C[A[r]] \leftarrow C[A[r]] - 1$ 
12  return  $B$ 
```

# El algoritmo CountingSort()

CountingSort ( $A, k$ ):

```
1    $B[0 \dots n-1] \leftarrow$  arreglo vacío
2    $C[0 \dots k] \leftarrow$  arreglo vacío
3   for  $i = 0 \dots k$  :
4        $C[i] \leftarrow 0$ 
5   for  $j = 0 \dots n-1$  :
6        $C[A[j]] \leftarrow C[A[j]] + 1$ 
7   for  $p = 1 \dots k$  :
8        $C[p] \leftarrow C[p] + C[p-1]$ 
9   for  $r = n-1 \dots 0$  :
10       $B[C[A[r]] - 1] \leftarrow A[r]$ 
11       $C[A[r]] \leftarrow C[A[r]] - 1$ 
12   return  $B$ 
```

- La complejidad del algoritmo es  $\Theta(n + k)$
- Si  $k \in \mathcal{O}(n)$ , entonces CountingSort() es  $\Theta(n)$

¡Este es un mejor tiempo que  $\mathcal{O}(n \log(n))$ !



## Ejemplo de ejecución

```
CountingSort (A, k):  
1    $B[0 \dots n-1] \leftarrow$  arreglo vacío  
2    $C[0 \dots k] \leftarrow$  arreglo vacío  
3   for  $i = 0 \dots k$  :  
4        $C[i] \leftarrow 0$   
5   for  $j = 0 \dots n-1$  :  
6        $C[A[j]] \leftarrow C[A[j]] + 1$   
7   for  $p = 1 \dots k$  :  
8        $C[p] \leftarrow C[p] + C[p-1]$   
9   for  $r = n-1 \dots 0$  :  
10       $B[C[A[r]] - 1] \leftarrow A[r]$   
11       $C[A[r]] \leftarrow C[A[r]] - 1$   
12   return  $B$ 
```

$A$

7	1	1	3	0	7	5	5	7	3
0	1	2	3	4	5	6	7	8	9

Hacemos el llamado CountingSort( $A, 7$ )

# Ejemplo de ejecución

CountingSort ( $A, k$ ):

- 1  $B[0 \dots n-1] \leftarrow$  arreglo vacío
- 2  $C[0 \dots k] \leftarrow$  arreglo vacío
- 3 **for**  $i = 0 \dots k$  :
- 4      $C[i] \leftarrow 0$

$A$

7	1	1	3	0	7	5	5	7	3
0	1	2	3	4	5	6	7	8	9

$C$

0	0	0	0	0	0	0	0
0	1	2	3	4	5	6	7

## Ejemplo de ejecución

CountingSort ( $A, k$ ):

```
1   $B[0 \dots n-1] \leftarrow$  arreglo vacío
2   $C[0 \dots k] \leftarrow$  arreglo vacío
3  for  $i = 0 \dots k$  :
4       $C[i] \leftarrow 0$ 
5  for  $j = 0 \dots n-1$  :
6       $C[A[j]] \leftarrow C[A[j]] + 1$ 
```

$A$

7	1	1	3	0	7	5	5	7	3
0	1	2	3	4	5	6	7	8	9

$C$

0	0	0	0	0	0	0	0
0	1	2	3	4	5	6	7

$C$

1	2	0	2	0	2	0	3
0	1	2	3	4	5	6	7

Hasta aquí,  $C[x]$  contiene el número de copias de  $x$  en  $A$

## Ejemplo de ejecución

CountingSort ( $A, k$ ):

```
1   $B[0 \dots n-1] \leftarrow$  arreglo vacío
2   $C[0 \dots k] \leftarrow$  arreglo vacío
3  for  $i = 0 \dots k$  :
4       $C[i] \leftarrow 0$ 
5  for  $j = 0 \dots n-1$  :
6       $C[A[j]] \leftarrow C[A[j]] + 1$ 
7  for  $p = 1 \dots k$  :
8       $C[p] \leftarrow C[p] + C[p-1]$ 
```

$A$

7	1	1	3	0	7	5	5	7	3
0	1	2	3	4	5	6	7	8	9

$C$

1	2	0	2	0	2	0	3
0	1	2	3	4	5	6	7

$C$

1	3	3	5	5	7	7	10
0	1	2	3	4	5	6	7

Hasta aquí,  $C[x]$  contiene cuántos elementos  
menores o iguales a  $x$  hay en  $A$

## Ejemplo de ejecución

```
9  for  $r = n - 1 \dots 0$  :  
10      $B[C[A[r]] - 1] \leftarrow A[r]$   
11      $C[A[r]] \leftarrow C[A[r]] - 1$ 
```

Para  $r = 9$

$A$

7	1	1	3	0	7	5	5	7	3
0	1	2	3	4	5	6	7	8	9

$B$

				3					
0	1	2	3	4	5	6	7	8	9

$C$

1	3	3	5	5	7	7	10
0	1	2	3	4	5	6	7

1	3	3	4	5	7	7	10
0	1	2	3	4	5	6	7

## Ejemplo de ejecución

Figure 1 illustrates the construction of the  $2^8$  Hadamard matrix  $H_8$  using the Kronecker product of the  $2^4$  Hadamard matrix  $H_4$  and the  $2^1$  Hadamard matrix  $H_2$ . The figure shows the recursive construction of  $H_8$  from  $H_4$  and  $H_2$ .

The  $2^4$  Hadamard matrix  $H_4$  is shown at the top, with rows labeled  $r=0$  to  $r=7$  and columns labeled  $c=0$  to  $c=3$ . The matrix is constructed from the  $2^2$  Hadamard matrix  $H_2$  and the  $2^1$  Hadamard matrix  $H_1$ .

The  $2^2$  Hadamard matrix  $H_2$  is shown in the middle, with rows labeled  $r=0$  to  $r=3$  and columns labeled  $c=0$  to  $c=1$ . The matrix is constructed from the  $2^1$  Hadamard matrix  $H_1$  and the  $2^0$  Hadamard matrix  $H_0$ .

The  $2^1$  Hadamard matrix  $H_1$  is shown at the bottom, with rows labeled  $r=0$  to  $r=1$  and columns labeled  $c=0$  to  $c=0$ . The matrix is constructed from the  $2^0$  Hadamard matrix  $H_0$  and the  $2^{-1}$  Hadamard matrix  $H_{-1}$ .

The figure demonstrates that the  $2^8$  Hadamard matrix  $H_8$  can be constructed as the Kronecker product of the  $2^4$  Hadamard matrix  $H_4$  and the  $2^1$  Hadamard matrix  $H_2$ , i.e.,  $H_8 = H_4 \otimes H_2$ .

# RadixSort

Otro algoritmo de ordenación en tiempo lineal es RadixSort

- Usado por las máquinas que ordenaban tarjetas perforadas
- Cada tarjeta tiene 80 columnas y 12 líneas
- Cada columna puede tener un agujero en una línea

El algoritmo ordena las tarjetas revisando una columna determinada

- Si hay un agujero en la columna, se pone en uno de los 12 compartimientos
- Las tarjetas con la perforación en la primera columna quedan encima
- La misma idea funciona para  $d$  columnas

Podemos generalizarla para un número natural de  $d$  dígitos

# Ordenando por dígito

Consideremos un número natural de  $d$  dígitos  $n_0 n_1 \dots n_{d-1}$

- Podemos ordenar según el dígito más significativo  $d_0$
- Según este dígito, separamos los números en *compartimientos*
- Luego, ordenados recursivamente cada compartimiento por su segundo dígito más significativo  $d_1$
- Finalmente, combinamos los contenidos de cada compartimiento

Problema: recursión busca que no mezclamos compartimientos antes de terminar



# Ordenación estable

Un ingrediente fundamental para el algoritmo que plantearemos es el siguiente

## Definición

Dada una secuencia  $A[0 \dots n-1]$ , sea  $B[0 \dots n-1]$  la secuencia resultante de ordenar  $A$  usando un algoritmo de ordenación  $S$ . Sean  $a, a'$  elementos en  $A$  tales que para el algoritmo  $A$  son equivalentes y  $a$  aparece antes que  $a'$  en  $A$ . Diremos que  $S$  es **estable** si los elementos correspondientes  $b$  y  $b'$  aparecen en el mismo orden relativo en  $B$ .

Si ordenamos por el segundo dígito, un orden estable dejaría elementos que comparten segundo dígito en el mismo orden en que nos llegaron

# RadixSort

El algoritmo RadixSort ordena por dígito **menos significativo**

- Ordena por dígito  $n_{d-1}$
- Luego, usando el mismo arreglo, ordena por dígito  $n_{d-2}$ , **con un algoritmo estable**
- Luego de ordenar  $k$  dígitos, los datos están ordenados si solo miramos el fragmento  $n_{d-k} \cdots n_{d-1}$
- Se requieren solo  $d$  pasadas para ordenar la secuencia completa

RadixSort( $A, d$ ):

**for**  $j = 0 \dots d - 1$  :

    StableSort( $A, j$ )   ▷ algoritmo de ordenación estable por  
       $j$ -ésimo dígito menos significativo

## Ejemplo de ejecución

	Arreglo inicial	Ordenado por unidad	Ordenado por decena	Ordenado por centena
0	0 6 4	0 0 <b>0</b>	0 <b>0</b> 0	<b>0</b> 0 0
1	0 0 8	0 0 <b>1</b>	0 <b>0</b> 1	<b>0</b> 0 1
2	2 1 6	5 1 <b>2</b>	0 <b>0</b> 8	<b>0</b> 0 8
3	5 1 2	3 4 <b>3</b>	5 <b>1</b> 2	<b>0</b> 2 7
4	0 2 7	0 6 <b>4</b>	2 <b>1</b> 6	<b>0</b> 6 4
5	7 2 9	1 2 <b>5</b>	1 <b>2</b> 5	<b>1</b> 2 5
6	0 0 0	2 1 <b>6</b>	0 <b>2</b> 7	<b>2</b> 1 6
7	0 0 1	0 2 <b>7</b>	7 <b>2</b> 9	<b>3</b> 4 3
8	3 4 3	0 0 <b>8</b>	3 <b>4</b> 3	<b>5</b> 1 2
9	1 2 5	7 2 <b>9</b>	0 <b>6</b> 4	<b>7</b> 2 9

# RadixSort

RadixSort( $A, d$ ):

**for**  $j = 0 \dots d - 1$  :

    StableSort( $A, j$ )   ▷ algoritmo de ordenación estable por  
       $j$ -ésimo dígito menos significativo

Supongamos que  $A$  tiene  $n$  datos naturales con  $d$  dígitos y se usa algoritmo estable lineal

- Si cada dígito puede tomar  $k$  valores distintos...
- Entonces RadixSort toma tiempo  $\Theta(d \cdot (n + k))$
- Si  $d$  es constante y  $k \in \mathcal{O}(n)$ , entonces RadixSort es  $\Theta(n)$

# Dos implementaciones

Estas ideas tiene dos implementaciones

LSD string sort (*Least Significant Digit*)

- Si todos los strings son del mismo largo (patentes, IP's, teléfonos)
- Funciona bien si el largo es pequeño

MSD string sort (*Most Significant Digit*)

- Si los strings tienen largo diferente
- Ordenamos con `CountingSort()` por primer caracter
- Recursivamente ordenamos subarreglos correspondientes a cada caracter (excluyendo el primero, que es común en cada subarreglo)
- Como Quicksort, puede ordenar de forma independiente
- **Pero** particiona en tantos grupos como valores del primer caracter

# MSD en acción

she  
sells  
seashells  
by  
the  
sea  
shore  
the  
shells  
she  
sells  
are  
surely  
seashells

are  
by  
she  
sells  
seashells  
sea  
shore  
shells  
she  
sells  
surely  
seashells  
the  
the

are  
by  
sells  
seashells  
sea  
sells  
seashells  
she  
shore  
shells  
she  
surely  
the  
the

are  
by  
seashells  
sea  
seashells  
sells  
sells  
she  
shore  
shells  
she  
surely  
the  
the

...

are  
by  
sea  
seashells  
seashells  
sells  
sells  
she  
she  
shells  
shore  
surely  
the  
the

# Cuidados de MSD

En la ejecución de MSD string sort se debe considerar

- Si un string  $s_1$  es prefijo de otro  $s_2$ ,  $s_1$  es menor que  $s_2$

$she \leq shells$

- Pueden usarse diferentes alfabetos

- binario (2)
- minúsculas (26)
- minúsculas + mayúsculas + dígitos (64)
- ASCII (128)
- Unicode (65.536)

- Para subarreglos pequeños (e.g.  $|A| \leq 10$ )

- cambiar a InsertionSort que *sepa* que los primeros  $k$  caracteres son iguales

# Sumario

Introducción

Funciones de hash

Ordenación lineal

**Cierre**



# Objetivos de la clase

- ☐ Identificar buenas propiedades de una función de hash
- ☐ Comprender algunas familias de funciones de hash comunes
- ☐ Conocer algoritmos de ordenación en tiempo lineal
- ☐ Comprender la limitación de dominio para tener tales algoritmos