

SESION 2: ÍNDICES

Ejercicio (Importante: índices multinivel secuencial, índice B+, índice Hash, coste de buscar con un índice)

Se dispone de un archivo r de datos que contiene información sobre estudiantes donde se almacena:

- Número de carnet: 8 bytes, todos los valores del campo son diferentes. El campo es un serial que comienza por valor 1. El archivo se encuentra ordenado por número de carnet
- Nombre del alumno: 40 bytes
- Código de carrera: 2 bytes con valores incluidos (IN) (0,1,2,3,4,5,6,7,8,9)
- Edad: 8 bytes
- Índice académico: 2 bytes

El número de registros a almacenar es de 100.000 y cada bloque tiene un tamaño de 512 bytes. El grado de ocupación de cada bloque es del 65%. La longitud de puntero a bloque es de 6 bytes y la longitud de puntero a registro de 7 bytes. Cada bloque tiene 12 bytes de control. Se pide para cada una de las situaciones que se muestran a continuación lo siguiente:

1. Calcular el número de bloques de un índice multinivel secuencial sobre el campo carnet de 3 niveles. Para las siguientes consultas determinar el número de registros estimados y el coste de lectura.

Clasificamos el índice \rightarrow primario + campo clave \rightarrow primer nivel puede ser:

- Denso $\rightarrow nri = V(\text{carnet})$
- Disperso $\rightarrow nri = Br$ (número de bloques de la tabla)

Si es denso $\rightarrow nri = V(\text{carnet}) = 100.000$ reg índice

$L_{ri} = L_{\text{carnet}} + L_{\text{preg}} = 8 + 7 = 15$ bytes.

$B_{\text{util}} = B - L_{\text{control}} = 512 - 12 = 500$ bytes * 0.65 = 325 bytes

$Fri = \lfloor 325 / 15 \rfloor = 21$ reg/bloq

$Bri = \lceil 100.000 / 21 \rceil = 4762$ bloq

Para segundo nivel y superior \rightarrow se hace índice disperso sobre el nivel anterior

$L_{ri2} = L_{\text{carnet}} + L_{\text{pbloque}} = 8 + 6 = 14$ bytes, $Fri = \lfloor 325 / 14 \rfloor = 23$ reg/bloq

$Bri2 = \lceil 4762 / 23 \rceil = 208$ bloq, nivel 2

$Bri3 = \lceil 208 / 23 \rceil = 10$ bloq, nivel 3

Número de bloques = $10 + 208 + 4762 = 4980$ bloq

Si es disperso $\rightarrow nri = br$ (una entrada por cada bloque de datos) = 20.000 reg índice

$L_{ri} = L_{\text{carnet}} + L_{\text{preg}} = 8 + 7 = 15$ bytes.

$$B_{\text{util}} = B - L_{\text{control}} = 521 - 12 = 500 \text{ bytes} * 0.65 = 325 \text{ bytes}$$

$$Fri = \lfloor 325 / 15 \rfloor = 21 \text{ reg/bloq}$$

$$Bri = \lceil 20.000 / 21 \rceil = 953 \text{ bloq}$$

Para segundo nivel y superior \rightarrow se hace índice disperso sobre el nivel anterior

$$Lri2 = L_{\text{carnet}} + L_{\text{pbloque}} = 8 + 6 = 14 \text{ bytes}, Fri = \lfloor 325 / 14 \rfloor = 23 \text{ reg/bloq}$$

$$Bri2 = \lceil 953 / 23 \rceil = 42 \text{ bloq, nivel 2}$$

$$Bri3 = \lceil 42 / 23 \rceil = 2 \text{ bloq, nivel 3}$$

$$\text{Número de bloques} = 2 + 42 + 953 = 997 \text{ bloq}$$

a. $\sigma_{\text{carnet}=2345}(r)$

$$Nrc = Nv * nr / V(\text{carnet}) = 1 * 100.000 / 100.000 = 1 \text{ reg}$$

Coste = Cíndice + Cdatos = Buscar nivel superior + 1 bloque por cada niveles inferiores al superior + bloques de la tabla a recuperar

$$\text{Para el denso: Coste} = \lceil \log_2(bri3) \rceil + 1 + 1 + \lceil nrc/fr \rceil = \lceil \log_2(10) \rceil + 1 + 1 + \lceil 1/5 \rceil = 4 + 1 + 1 + 1 = 7 \text{ bloq}$$

$$\text{Para el disperso: Coste} = \lceil \log_2(bri3) \rceil + 1 + 1 + \lceil nrc/fr \rceil = \lceil \log_2(2) \rceil + 1 + 1 + \lceil 1/5 \rceil = 1 + 1 + 1 + 1 = 4 \text{ bloq}$$

$$\text{Búsqueda secuencial} = \text{Coste} = 20.000 \text{ blq (peor)}$$

$$\text{Búsqueda binaria} = \lceil \log_2(20.000) \rceil + \lceil 1/5 \rceil - 1 = 15 \text{ blq}$$

b. $\sigma_{\text{código_carrera}=2}(r)$

$$Nrc = Nv * nr / V(\text{código_carnet}) = 1 * 100.000 / 10 = 10.000 \text{ blq}$$

$$\text{Solo búsqueda secuencial} \rightarrow \text{Coste} = br = 20.000 \text{ blq}$$

2. Calcular el número de bloques de un índice multinivel secuencial sobre el campo código_carrera. Para las siguientes consultas determinar el número de registros estimados y el coste de lectura.

Clasificamos el índice \rightarrow secundario + campo no clave \rightarrow cajones de punteros

$$\text{Primer nivel solo denso} \rightarrow nri = V(\text{código_carrera}) = 10$$

$$\text{Como es denso} \rightarrow nri = V(\text{código_carnet}) = 10 \text{ reg índice}$$

$$Lri = L_{\text{código_carrera}} + L_{\text{bloques}} = 2 + 6 = 8 \text{ bytes.}$$

$$B_{\text{util}} = B - L_{\text{control}} = 521 - 12 = 500 \text{ bytes} * 0.65 = 325 \text{ bytes}$$

$$Fri = \lfloor 325 / 8 \rfloor = 40 \text{ reg/bloq}$$

$$Bri = \lceil 10 / 40 \rceil = 1 \text{ bloq}$$

Como solo hay 1 bloque en el nivel inferior \rightarrow solo tiene un nivel

$$\text{Cajones de punteros} \rightarrow \text{Número de cajones } Nc = V(\text{código_carnet}) = 10$$

$$\text{Cada cajon tiene un número de registros } nc = nr / V(\text{código_carnet}) = 100.000 / 10 = 10.000 \text{ reg}$$

Los cajones de punteros llevan los Pregistro \rightarrow $Lrc = LPreg = 7$ bytes \rightarrow
 $Frc = \lfloor 325 / 7 \rfloor = 46$ reg/bloq \rightarrow $brc = \lceil 10.000 / 46 \rceil = 218$ blq
 Número de bloques del índice = $bri + Nc * brc = 1 + 10 * 218 = 2181$ bloq

a. $\sigma_{carnet=2345}(r)$

Solo búsqueda secuencial \rightarrow Coste = $br = 20.000$ bloq

b. $\sigma_{código_carrera=2}(r)$

$Nrc = Nv * nr / V(código_Carnet) = 1 * 100.000 / 10 = 10.000$ reg (ordenados
 \rightarrow Consecutivos)

Coste = $1 + 1$ cajon de punteror + $nrc = 1 + 218 + 10.000 = 10.219$ bloq

Búsqueda secuencial \rightarrow Coste = $br = 20.000$ bloq

3. Calcular el número de bloques de un índice B+ sobre el campo carnet. Para cada una de las siguientes consultas determinar el número de registros estimados y el coste de lectura.

Clasificamos el índice \rightarrow primario + campo clave.

Los B+ siempre densos \rightarrow $nri = V(carnet) = 100.000$ reg

Nodo raíz/intermedio

$n * Lpbloque + (n-1) * Lcarnet \leq Butil \rightarrow n * 6 + (n-1) * 8 \leq 325 \rightarrow 14 * n \leq 333 \rightarrow n = 23$ punteros a bloque

Nodo hoja

$nh * (Lcarnet + Lpreg) + Lpbloque \leq Butil \rightarrow nh * (8 + 7) + 6 \leq 325 \rightarrow 15 * nh \leq 319 \rightarrow n = 21$ valores del campo

Nodos hoja = $\lceil 100.000 / 21 \rceil = 4762$ bloq

Nodos intermedio 1 = $\lceil 4762 / 23 \rceil = 208$ bloq

Nodos intermedio 2 = $\lceil 208 / 23 \rceil = 9$ bloq

Nodo raíz = $\lceil 9 / 23 \rceil = 1$ bloq

Número de bloques = $1 + 9 + 208 + 4762 = 4980$ bloq, Son 4 niveles

a. $\sigma_{carnet=2345}(r)$

$Nrc = 1$ reg, como vimos antes.

Para buscar los registros:

Número Búsqueda secuencial \rightarrow Coste = $br = 20.000$ bloq

Búsqueda en el B+ (1 sola) = niveles + $nrc / fr = 4 + 1 = 5$ bloq

b. $\sigma_{código_carrera=2}(r)$

$Nrc = 10.000$ reg, como vimos antes.

Para buscar los registros:

Búsqueda secuencial \rightarrow Coste = br = 20.000 bloq

4. Calcular el número de bloques de un índice B+ sobre el campo código_carrera.
Para cada una de las siguientes consultas determinar el número de registros estimados y el coste de lectura.

Clasificamos el índice \rightarrow secundario + campo no clave \rightarrow Cajones de punteros

Los B+ siempre densos \rightarrow nri = V(código_carrera) = 10 reg

Nodo raíz/intermedio

$$n * L_{\text{pbloque}} + (n-1) * L_{\text{codigo_carrear}} \leq \text{Butil} \rightarrow n * 6 + (n-1) * 2 \leq 325 \rightarrow 8 * n \leq 327 \rightarrow n = 40 \text{ punteros a bloque}$$

Nodo hoja

$$nh * (L_{\text{codigo_carrera}} + L_{\text{bloque}}) + L_{\text{pbloque}} \leq \text{Butil} \rightarrow nh * (2+6) + 6 \leq 325 \rightarrow 8 * nh \leq 319 \rightarrow n = 39 \text{ valores del campo}$$

$$\text{Nodos hoja/Raíz} = \lceil 10 / 39 \rceil = 1 \text{ bloq}$$

Número de bloques = 1 bloq, es 1 nivel

Hay cajones de punteros. Lo mismo que en el apartado 2.

Cajones de punteros \rightarrow Número de cajones $N_c = V(\text{código_carnet}) = 10$

Cada cajon tiene un número de registros
 $nc = nr / V(\text{código_carnet}) = 100.000 / 10 = 10.000 \text{ reg}$

Los cajones de punteros llevan los Pregiro \rightarrow Lrc = LPreg = 7 bytes \rightarrow
 $Frc = \lfloor 325 / 7 \rfloor = 46 \text{ reg/bloq} \rightarrow brc = \lceil 10.000 / 46 \rceil = 218 \text{ blq}$

a. $\sigma_{\text{carnet}=2345}(r)$

Nrc = 1 reg, como vimos antes.

Para buscar los registros:

Búsqueda secuencial \rightarrow Coste = br = 20.000 bloq

b. $\sigma_{\text{código_carrera}=2}(r)$

Nrc = 10.000 reg, como vimos antes.

Para buscar los registros:

Búsqueda secuencial \rightarrow Coste = br = 20.000 bloq

Búsqueda en el B+ (1 sola) = Número niveles + 1 cajón de punteros + los registros a los que apunta el cajón (nrc) = 1 + 1*218+10.000 = 10.219 blq

5. Calcular el número de bloques de un índice Hash creado sobre el campo carnet cuya **función de asociación devuelve un número binario de 8 bits**.

Para cada una de las siguientes consultas determinar el número de registros estimados y el coste de lectura.

Clasificamos el índice \rightarrow primario + campo clave

Los hash son siempre densos $\rightarrow nri = V(\text{campo}) = nr = 100.000 \text{ reg}$

Se examina la función de asociación \rightarrow número binario de 8 bits \rightarrow Número de cajones $\rightarrow N = 2^8 = 256 \text{ cajones}$.

Cada cajón lleva de media $nc = nri / N = 100.000 / 256 = 390,625 \text{ reg} \rightarrow 391 \text{ reg}$

$Lric = Lcarnet + Lpreg = 8 + 7 = 15 \text{ reg}$, $frc = \lfloor 325 / 15 \rfloor = 21 \text{ reg / bloq}$, $bc = \lceil 391 / 21 \rceil = 19 \text{ bloq}$

El índice ocupa : $Bri = N * Bc = 256 * 19 = 4864 \text{ blq}$

a. $\sigma_{\text{carnet}=2345}(r)$

$Nrc = 1 \text{ reg}$, como vimos antes.

Para buscar los registros:

Búsqueda secuencial \rightarrow Coste = $br = 20.000 \text{ bloq}$

Búsqueda binaria = $\lceil \log_2 (20.000) \rceil + \lceil 1/5 \rceil - 1 = 15 \text{ blq}$

Búsqueda por Índice Hash \rightarrow 1 valor en 1 cajón + $\lceil nrc/fr \rceil = 19 + 1 = 20 \text{ blq}$

b. $\sigma_{\text{código_carrera}=2}(r)$

$Nrc = 10.000 \text{ reg}$, como vimos antes.

Para buscar los registros:

Búsqueda secuencial \rightarrow Coste = $br = 20.000 \text{ bloq}$

6. Calcular el número de bloques de un índice Hash creado sobre el campo `código_carrera` cuya función de asociación es **`código_carrera mod 5`**. Para cada una de las siguientes consultas determinar el número de registros estimados y el coste de lectura.

Clasificamos el índice \rightarrow secundario + campo no clave \rightarrow cajones de punteros

Los hash son siempre densos $\rightarrow nri = V(\text{campo}) = nr = 10 \text{ reg}$

Se examina la función de asociación \rightarrow `código_carrera mod 5` \rightarrow Número de cajones $\rightarrow N = 5$ (del 0 al 4)

Cada cajón lleva de media $nc = nri / N = 10 / 5 = 2 \text{ reg}$

$Lric = Lcarnet + Lbloque = 2 + 6 = 8 \text{ reg}$, $frc = \lfloor 325 / 8 \rfloor = 40 \text{ reg / bloq}$, $bc = \lceil 2 / 40 \rceil = 1 \text{ bloq}$

Hay cajones de punteros. Lo mismo que en el apartado 2.

Cajones de punteros \rightarrow Número de cajones $Nc = V(\text{código_carnet}) = 10$

Cada cajon tiene un número de registros
 $nc = nr / V(\text{código_carnet}) = 100.000 / 10 = 10.000 \text{ reg}$

Los cajones de punteros llevan los Pregistro $\rightarrow Lrc = LPreg = 7 \text{ bytes} \rightarrow$
 $Frc = \lfloor 325 / 7 \rfloor = 46 \text{ reg/bloq} \rightarrow brc = \lceil 10.000 / 46 \rceil = 218 \text{ blq}$

El índice ocupa : $Bri = N * Bc + V(\text{código_Carrera}) * Brc = 5 * 1 + 10 * 218 = 2185 \text{ blq}$

a. $\sigma_{\text{carnet}=2345}(r)$

$Nrc = 1 \text{ reg}$, como vimos antes.

Para buscar los registros:

Búsqueda secuencial $\rightarrow \text{Coste} = br = 20.000 \text{ bloq}$

Búsqueda binaria $= \lceil \log_2 (20.000) \rceil + \lceil 1/5 \rceil - 1 = 15 \text{ blq}$

b. $\sigma_{\text{código_carrera}=2}(r)$

$Nrc = 10.000 \text{ reg}$, como vimos antes.

Para buscar los registros:

Búsqueda secuencial $\rightarrow \text{Coste} = br = 20.000 \text{ bloq}$

Búsqueda índice hash $\rightarrow 1 \text{ valor a buscar en } 1 \text{ cajón} + \text{ leer el cajón de punteros correspondiente} + \text{ los registros que apuntan } nrc = 1 + 218 + 10.000 = 10.219 \text{ bloq}$

7. Cuestiones opcionales para resolver por los alumnos. Considerando los tipos de índices del ejercicio anterior, determinar el número de registros a recuperar y el coste de buscar para cada una de las siguientes consultas.

a. $\sigma_{\text{carnet} \geq 10.000}(r)$

Número de registros a recuperar $nrc = 90.001 * 100.000 / 100.00 = 90.001 \text{ reg}$

Cómo se puede buscar:

- Secuencial $\rightarrow \text{Coste} = 20.000 \text{ bloq}$
- Binaria $\rightarrow \text{Coste} = \lceil \log_2 (20.000) \rceil + \lceil 90.001/5 \rceil - 1 = 15 + 18.001 - 1 = 18.015 \text{ bloq}$
- B+ en carnet, como es primario, se busca valor 10.000 luego se leen 90.001 registros consecutivos $\rightarrow \text{Coste} = 4 + \lceil 90.001/5 \rceil = 18.004 \text{ bloq}$
- Hash en carnet, como es primario, se busca valor 10.000 y luego se leen 90.001 registros consecutivos $\rightarrow \text{Coste} = 19 (1 \text{ Cajón}) + \lceil 90.001/5 \rceil = 18.019 \text{ bloq}$
-

b. $\sigma_{\text{carnet} < 10.000}(r)$

Número de registros a recuperar $nrc = 9999 * 100.000 / 100.00 = 9.999$ reg

Cómo se puede buscar:

- Secuencial \rightarrow Coste = 20.000 bloq
- Secuencial limitada \rightarrow Solo hay que leer 9.999 reg empezando por el principio del archivo ordenado \rightarrow Coste = $\lceil 9.999/5 \rceil = 2000$ bloq
- Binaria \rightarrow Coste = $\lceil \log_2(20.000) \rceil + \lceil 9.999/5 \rceil - 1 = 15 + 2000 - 1 = 2014$ bloq
- B+ en carnet, como es primario, se busca valor 1 y luego se leen 9.999 registros consecutivos \rightarrow Coste = $4 + \lceil 9.999/5 \rceil = 2004$ bloq
- Hash en carnet, como es primario, se busca valor 1 y luego se leen 9.999 registros consecutivos \rightarrow Coste = 19 (1 Cajón) + $\lceil 9999/5 \rceil = 2019$ bloq

c. $\sigma_{\text{carnet} < > 10.000}(r)$

Número de registros a recuperar $nrc = 99999 * 100.000 / 100.00 = 99.999$ reg. Queremos todos menos el 10.000

Cómo se puede buscar:

- Secuencial \rightarrow Coste = 20.000 bloq
- Binaria \rightarrow No están todos los registros seguidos, luego leer desde el 1 hasta al 9.999, y luego localizar el 10001 para leer hasta el 100.000. En este caso al ser un campo con todos los valores diferentes y al querer todos menos 1, mejor secuencial.
- B+ en carnet, ocurre lo mismo que antes.
- Hash en carnet, ocurre lo mismo que antes, hay que leerse todos los bloques del archivo, luego saldrá más coste que la secuencial.

d. $\sigma_{\text{carnet} > 10.000 \wedge \text{carnet} < 15.000}(r)$

Número de registros a recuperar $nrc = (14999 - 10001 + 1) * 100.000 / 100.000 = 4.999$ reg. Desde el 10.001 al 14.999, incluidos.

Cómo se puede buscar:

- Secuencial \rightarrow Coste = 20.000 bloq
- Binaria \rightarrow Coste = $\lceil \log_2(20.000) \rceil + \lceil 4.999/5 \rceil - 1 = 15 + 1000 - 1 = 1014$ bloq. Busco el primer valor que cumple y luego leo los registros consecutivamente del archivo de datos que para eso está ordenado.

- B+ en carnet \rightarrow Coste = $4 + \lceil 4.999/5 \rceil = 4 + 1000 = 1004$ bloq. Localizo el primer valor que cumple la condición con el árbol y luego tiro por el archivo de datos.
- Hash en carnet \rightarrow Coste = 19 (1 Cajón) + $\lceil 4.999/5 \rceil =$ Localizo el primer valor que cumple condición y luego sigo por el archivo de datos.

e. $\sigma_{\text{código_carrera} > 6} (r)$

Número de registros a recuperar $nrc = (9-7+1)*100.000/10 = 30.000$ reg. Se pide el 7,8,9. El archivo no está ordenado, luego sacar todos los registros a partir del índice.

Cómo se puede buscar:

- Secuencial \rightarrow Coste = 20.000 bloq
- B+ en carnet \rightarrow Coste = 1 (nodo raíz/hoja para leer el puntero a cajones de punteros a registro del 7,8,9) + $3*218 + 3*10.000 = 30.655$ bloq. Peor que la secuencial.
- Hash en carnet \rightarrow Coste = $3*1$ (caen en 3 cajones diferentes) + $3*218 + 3*10.000 = 30.657$ bloq. Peor que la secuencial.

f. $\sigma_{\text{código_carrera} < 5} (r)$

Número de registros a recuperar $nrc = (4-0+1)*100.000/10 = 50.000$ reg. Se pide el 0,1,2,3,4. El archivo no está ordenado, luego sacar todos los registros a partir del índice.

Cómo se puede buscar:

- Secuencial \rightarrow Coste = 20.000 bloq
- B+ en carnet \rightarrow Coste = 1 (nodo raíz/hoja para leer el puntero a cajones de punteros a registro del 0,1,2,3,4) + $5*218 + 5*10.000 = 51.091$ bloq. Peor que la secuencial.
- Hash en carnet \rightarrow Coste = $5*1$ (caen en 5 cajones diferentes) + $5*218 + 5*10.000 = 51.095$ bloq. Peor que la secuencial.

g. $\sigma_{\text{código_carrera} < > 8} (r)$

Número de registros a recuperar $nrc = (10-1)*100.000/10 = 90.000$ reg. Se pide todos menos el 8. El archivo no está ordenado, luego sacar todos los registros a partir del índice.

Cómo se puede buscar:

- Secuencial \rightarrow Coste = 20.000 bloq
- B+ en carnet \rightarrow Coste = 1 (nodo raíz/hoja para leer el puntero a cajones de punteros a registro del 0,1,2,3,4,5,6,7,9) + $9*218 + 9*10.000 = 91.936$ bloq. Peor que la secuencial.

- Hash en carnet \rightarrow Coste = $5*1$ (caen en 5 cajones diferentes) + $9*218 + 9*10.000 = 91.940$ bloq. Peor que la secuencial.

h. $\sigma_{\text{código_carrera} > 7 \wedge \text{código_carrera} \leq 9}(r)$

Número de registros a recuperar $nrc = (9-8+1)*100.000/10 = 20.000$ reg. Se pide el 8,9. El archivo no está ordenado, luego sacar todos los registros a partir del índice.

Cómo se puede buscar:

- Secuencial \rightarrow Coste = 20.000 bloq
- B+ en carnet \rightarrow Coste = 1 (nodo raíz/hoja para leer el puntero a cajones de punteros a registro del 8,9) + $2*218 + 2*10.000 = 20.437$ bloq. Peor que la secuencial.
- Hash en carnet \rightarrow Coste = $2*1$ (caen en 2 cajones diferentes) + $2*218 + 2*10.000 = 20.438$ bloq. Peor que la secuencial.

i. $\sigma_{\text{carnet} = 53456 \wedge \text{código_carrera} = 9}(r)$

Número de registros a recuperar $nrc = 100.000/(100.000*10) = 0.1$ reg. Será o ningún registro o 1 como máximo debido al campo carnet.

Como es un AND significa que la tabla se puede buscar por uno de ellos, el más barato, que como se ha visto antes es por carnet. Luego el coste sería como el visto en los apartados anteriores con condición de carnet = 2345. Si se busca por código de carrera al devolver más registros entonces su coste es mayor y además el archivo no está ordenado por ee campo. En este caso para Carnet el B+ da menor coste.

j. $\sigma_{\text{carnet} = 53456 \vee \text{código_carrera} = 9}(r)$

Número de registros a recuperar $nrc = 100.000/100.000 + 10.000/10 = 10.001$ reg. Es un OR sobre dos campos luego equivale a realizar dos búsquedas una por cada campo si no se usa la lectura secuencial.

Coste será:

- \rightarrow Coste de carnet=53456 que si se tienes un B+ es más barato para este problema con Coste = $4+1 = 5$ bloques. Obtendría 1 reg o 0 registros.
- \rightarrow Coste de código de carrera=9, El hash o el B+ en código de carrera porque su coste es: Coste = $1 + 218 + 10.000 = 10.219$ bloq en este caso.
- \rightarrow La secuencial tiene un coste de 20.000 bloques, luego es más barato buscar por los dos índices: B+ por carnet y B+ o Hash por código de carrera teniendo un coste total de $5 + 1 + 218 + 10000 = 10.224$ bloq