

Pauta

Control 2 - Indices y algoritmos

El plazo para la entrega es el viernes 14 de octubre, hasta las 23:59. La entrega debe hacerse en un **pdf** y subirse a Canvas (se creará una tarea con el form para subir la entrega).

Pregunta 1 (2 pts)

Considere la relación

Cervezas(id int PRIMARY KEY, nombre varchar(20), marca varchar(20), estilo varchar(15), precio int)

indexada con un hash-join clustered sobre el atributo id, un hash-join unclustered sobre marca y un B+tree unclustered sobre precio.

1. Indique el costo I/O de las siguientes consultas asumiendo que la relación tiene 10000 tuplas, que el id va desde el 0 hasta el 9999, que en una pagina caben 80 tuplas o 250 punteros, que la altura del B+Tree es h , que sus hojas están ocupadas a un 60% y que el hash index tiene solo una pagina por bucket.

- SELECT * FROM Cervezas WHERE id = 20
- SELECT * FROM Cervezas WHERE id <=3
- SELECT * FROM Cervezas WHERE precio >2000 AND precio <3000 (asuma que existen menos de 80 tuplas que cumplen esta condición)
- SELECT * FROM Cervezas WHERE marca = 'Tamango Brebajes'
- SELECT * FROM Cervezas WHERE estilo = 'Hazy IPA'

0.3 c/u

2. Bajo las mismas condiciones de la pregunta anterior, explique como podría resolver la consulta SELECT * FROM Cervezas WHERE precio <3000 AND marca = 'Tamango Brebajes' (haciendo un uso inteligente de los indices) e indique el costo I/O de su propuesta. Asuma que la cantidad de cervezas que valen menos de 3000 caben en una página. → 0.5.

Pregunta 2 (2 pts)

Considere que tiene las relaciones $R(a \text{ int PRIMARY KEY}, b \text{ int})$ y $S(b \text{ int PRIMARY KEY}, c \text{ int})$. Escriba en pseudo-código un programa para computar $R \bowtie S$ tal que recorra R y S solo una vez. Hint: es probable que tenga que pre-procesar la(s) relacion(es) de alguna manera...

Pregunta 3 (2 pts)

Considere las relaciones $R(a \text{ int}, b \text{ int})$ y $S(a \text{ int}, c \text{ int})$ e imagine que casi todas las consultas que se hacen implican hacer un join entre R y S . Imagine también que R ya se encuentra indexada según el atributo b . Responda lo siguiente:

1. Explique cómo resolver $R \bowtie S$ y cuál es su costo en operaciones I/O 0.5
2. Recomendaría usted incluir algún otro índice en las relaciones? Si la respuesta es sí, indique el tipo de índice y cómo quedaría el costo I/O de la consulta. Si la respuesta es no, indique por qué no. 1.0.
3. ¿Existe alguna manera en la que se puede disminuir el costo I/O de la consulta sin crear nuevos índices? Justifique su respuesta. 0.5

Pregunta #1

1. Indique el costo I/O de las siguientes consultas asumiendo que la relación tiene 10000 tuplas, que el id va desde el 0 hasta el 9999, que en una pagina caben 80 tuplas o 250 punteros, que la altura del B+Tree es h , que sus hojas están ocupadas a un 60% y que el hash index tiene solo una pagina por bucket. (1.5 pts)

- SELECT * FROM Cervezas WHERE id = 20
- SELECT * FROM Cervezas WHERE id <= 3
- SELECT * FROM Cervezas WHERE precio > 2000 AND precio < 3000 (asuma que existen menos de 80 tuplas que cumplen esta condición)
- SELECT * FROM Cervezas WHERE marca = 'Tamango Brebajes'
- SELECT * FROM Cervezas WHERE estilo = 'Hazy IPA'

(0.3 pts

cada una)

• SELECT * FROM Cervezas WHERE id = 20

El costo es $O(1)$ pues usamos el hash-index

• SELECT * FROM Cervezas WHERE id <= 3

El costo es 4, pues en el pen con cada tupla estara en un bucket distinto.

• -. precio > 2000 AND precio < 3000
Aquí usamos el B+tree, el costo es

$\frac{h}{\downarrow}$	$+$	$\frac{1}{\downarrow}$	$+$	$\frac{x}{\downarrow}$	
altura del árbol		pagina donde estan los punteros		tuplas que cumplen la condición (son menos de 80)	

↳ podrian acotarlo pm • $h + 81$
• $h + 80$

• ... marca = 'Tamango Brebajes'

↓
marca está indexada en un hash-index
en clustered, entonces el costo I/O es

$$1 + X$$

↳ cantidad de tuplas que
cumplen la condición

• ... estilo = "Hazy IPA"

↓
como no hay índice sobre estilo, hay que
recurrir a todas las páginas

$$\Rightarrow \text{el costo es } \frac{10.000}{80} = 125$$

2. Bajo las mismas condiciones de la pregunta anterior, explique como podría resolver la consulta `SELECT * FROM Cervezas WHERE precio < 3000 AND marca = 'Tamango Brebajes'` (haciendo un uso inteligente de los índices) e indique el costo I/O de su propuesta. Asuma que la cantidad de cervezas que valen menos de 3000 caben en una página. (0.5 pts)

↳ Esto es lo mismo que la penúltima
consulta, busco la marca en el
hash-index clustered y luego voy
cada puntero ...

la única dif. es que tengo que revisar
una condición más, pero eso no afecta
el costo I/O

$$\text{Costo I/O} \quad 1 + X$$

* Oña opam de la pregunta 1.2 sea sea
el R tal en memo, en cuyo caso el
costo sea $R + X$
lo # de veces en
memo < 3000 .

Pregunta 2

Considere que tiene las relaciones $R(a \text{ int PRIMARY KEY}, b \text{ int})$ y $S(b \text{ int PRIMARY KEY}, c \text{ int})$. Escriba en pseudo-código un programa para computar $R \bowtie S$ tal que recorra R y S solo una vez. *Hint: es probable que tenga que pre-procesar la(s) relacion(es) de alguna manera...* 2 pts.

La idea es simple: primero hay que ordenar las relaciones y después hacer lo siguiente

```
R.open()
S.open()
x = R.next()
y = S.next()
while x != NULL AND y != NULL:
    if  $\pi_b(x) = \pi_b(y)$ :
        OUTPUT x
        x = R.next()
        y = S.next()
    elif x < y:
        x = R.next()
    elif y < x:
        y = S.next()
```

Como están
ordenados, solo
avanzar en el
más pequeño
y así se llama
sort merge
join

Pregunta 3

Considere las relaciones $R(a \text{ int}, b \text{ int})$ y $S(a \text{ int}, c \text{ int})$ e imagine que casi todas las consultas que se hacen implican hacer un join entre R y S . Imagine también que R ya se encuentra indexada según el atributo b . Responda lo siguiente:

1. Explique cómo resolver $R \bowtie S$ y cuál es su costo en operaciones I/O **0.5 pts**
2. Recomendaría usted incluir algún otro índice en las relaciones? Si la respuesta es sí, indique el tipo de índice y cómo quedaría el costo I/O de la consulta. Si la respuesta es no, indique por qué no. **1 pt**
3. ¿Existe alguna manera en la que se puede disminuir el costo I/O de la consulta sin crear nuevos índices? Justifique su respuesta. **0.5 pts**

3.1) Podemos hacer un nested loop join
Cuyo costo es

$$\text{Costo}(R) + \text{tuples}(R) \cdot \text{Costo}(S)$$

→ acá los estudiantes pueden proponer tb
un block nested loop join.

3.2) lo más lógico sería indexar R , S
o ambos según el atributo ~~b~~.a. Con el
costo de R el índice sería en clustered,
en el caso de S , como no hay otro
índice, podría ser clustered.

Si usamos hash index en ambos,
el costo podría ser

$$\text{Costo}(R) + \text{tuples}(R)$$

$$\text{Costo}(S) + 2 \cdot \text{tuples}(S)$$

Si usas un
B+ tree podría
ser distinto...

3.3) Una manera de agilizar la consulta sin nuevos índices sería ordenar las relaciones...