

Parcialito 5 - Optimización de Consultas - CORRECCIÓN RÁPIDA

● Graded

Student

Juana Rehl

Total Points

10 / 10 pts

Question 1

Estime el costo de realizar la consulta y la cantidad de filas que serán devueltas.

10 / 10 pts

✓ - 0 pts Correcto

- 1 pt Mal calculada la cardinalidad por fechas
- 1 pt Error de cálculos
- 2.5 pts Falta calcular cantidad de filas de la salida (cardinalidad)
- 2 pts Error conceptual
- 2 pts Falta probar casos de join

1 Por qué no aparecería $B(R)$ o $B(R')$?

2 Aparte no tenés índice para hacer la junta después

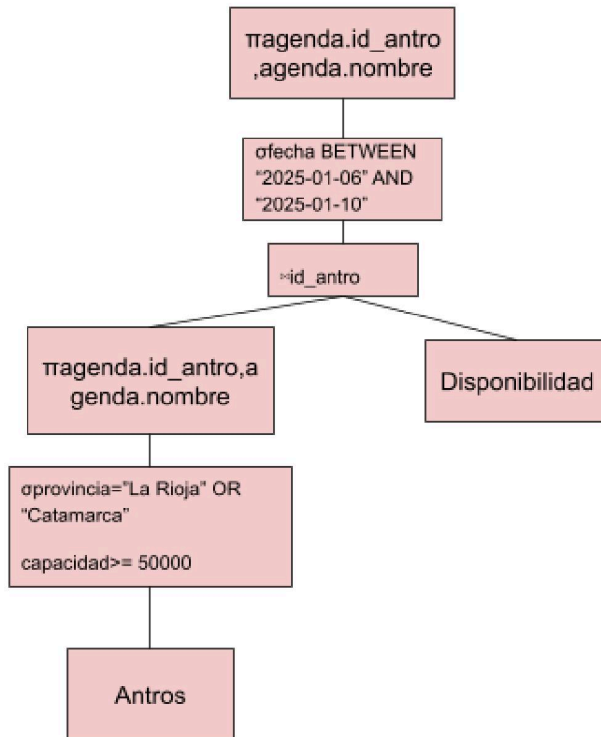
Question assigned to the following page: [1](#)

Juana Rehl - 112185 - Parcialito 5

Resolución

R = Antros

S = Disponibilidad



1. Aplicación de selección sobre R. Primero se filtra la relación Antros por provincia y capacidad:

$$R' = \sigma((\text{provincia} = \text{'La Rioja'} \text{ OR } \text{provincia} = \text{'Catamarca'}) \wedge (\text{capacidad} \geq 50000))(R)$$

Costo por capacidad:

- FileScan $\rightarrow B(R) = 2000$ bloques, (No hay índice sobre capacidad).

Costo por provincia:

- IndexScan $\rightarrow H(I(A, R)) + \lceil n(R)/V(A, R) \rceil = 2 + 10000/20 = 502$.
Como se trata de dos valores (por el OR), se duplica: $502 \times 2 = 1004$.

Se elige el método por provincia, ya que resulta ser el de menor costo en comparación al de capacidad. $\text{COST}(\sigma_{\text{provincia}, \text{capacidad}}(R)) = 1004$.

Question assigned to the following page: [1](#)

Cantidad de filas filtradas:

$$n(R') = (10000 \times 2 / 20) \times 0.1 = 100$$

Para poder sacar el $n(R')$ se tiene en cuenta que de las 20 provincias se precisan 2 de ellas y tenemos información de que el 10% de los antros tienen capacidad para al menos 50 mil personas (por eso lo multiplicamos por 0.1).

$$\text{Factor de bloque: } F(R) = 10000 / 2000 = 5$$

$$\text{Cantidad de bloques: } B(R') = 100 / 5 = 20$$

En este punto, se podría plantear otra alternativa e ir por otro enfoque. Por ejemplo en vez de seleccionar por provincias/capacidad, seleccionar según la fecha (ofecha(s)). Aunque es una alternativa completamente válida, rápidamente nos podemos dar cuenta que no es la más óptima.

Si realizamos un FileScan por fecha: FileScan \rightarrow 5000. Por IndexScan $\rightarrow 3 + 100000 * 5/100 = 5003$. Entre estas dos opciones nos conviene el FileScan, pero si lo comparamos con el enfoque seleccionado anteriormente, nos damos cuenta que el costo de seleccionar por provincias/capacidad es de 1004 que es mucho menor que 5000. Descartamos esta alternativa por fechas.

2. Evaluación del JOIN ($R' \bowtie S$). Disponemos de $M = 100$ bloques de memoria.

a) Loops anidados por bloque:

$$\text{COST} = \lceil B(R') / (M-2) \rceil \times B(S) = 1 \times 5000 = 5000$$

b) Loop con único índice:

$$\text{COST} = n(R') \times (H(I(A, S)) + \lceil n(R') / V(A, S) \rceil) = 100 \times (4 + 100000/10000) = 1400$$

c) Sort-Merge:

$$\text{COST} = B(S) + 2 \times B(R') \times \lceil \log_{M-1}(B(R')) \rceil + 2 \times B(S) \times \lceil \log_{M-1}(B(S)) \rceil = 10000 + 2 \times 20 \times 1 + 2 \times 10000 \times 1 = 30040$$

d) Hash Join:

Límites: $P \leq M - 1 \leq 99$ y $\min(1; 51) \leq M - 2 \leq 98$

$$\text{COST} = 2 \times B(R') + 3 \times B(S) = 2 \times 20 + 3 \times 5000 = 15040$$

El método de loop con único índice resulta el más eficiente, con costo total = 1400.

Cantidad de filas luego del join:

$$n(R' \bowtie S) = 100 \times (100000 / 10000) = 1000$$

Factor de bloque:

$$F(R' \bowtie S) = 5 \times 20 / (5 + 20) = 4$$

$$\text{Cantidad de bloques: } B(R' \bowtie S) = 1000 / 4 = 250$$

Question assigned to the following page: [1](#)

3. Filtro por fecha. Sea $Y = R \bowtie S$. Aplicamos el filtro sobre fecha (entre el 6 y el 10 de enero). El operador se ejecuta en pipeline, por lo que no incurre en costo adicional de lectura. $COST(\sigma_{\text{fecha}}(Y)) = 0$ ✓

4. Proyección final

Como no se requiere eliminar duplicados, la proyección se realiza directamente sobre la salida del operador anterior, sin costo extra. ✓

Costo total estimado

$$C_{\text{total}} = 1004 \text{ (selección)} + 1400 \text{ (join)} = 2404$$

Filas estimadas

$$n(Y) = 1000 \times 0.05 = 50 \text{ (ya que se toman 5 de las 100 fechas posibles)} \quad \checkmark$$

El plan más eficiente consiste en usar el join con índice único, con un costo total aproximado de 2404 bloques y una salida esperada de 50 registros.