

Parcialito 6

Alumno: Alan Ezequiel Valdevenito.
Padrón: 107585.

Ejercicio 1

Item a

Esquema
profesores(<u>legajo</u> , nombre, apellido, genero, titulo)

Consulta
$\sigma_{titulo="LICENCIADO" \wedge genero="M"}(profesores)$

Información
<ul style="list-style-type: none">▪ $n(\text{profesores}) = 500$ y $F(\text{profesores}) = 25$.▪ $V(\text{titulo}, \text{profesores}) = 10$ y $V(\text{genero}, \text{profesores}) = 2$.▪ La tabla profesores tiene un índice llamado prtít por título con $\text{Height}(\text{prtít}) = 2$. Este índice no es de clustering.▪ Los valores de las columnas titulo y genero se almacenan siempre en mayúsculas.

Esta consulta tiene como resultado aquellas filas de la tabla profesores tales que se cumpla simultáneamente la condición de que su titulo sea LICENCIADO y su género sea M.

Notemos la tabla profesores tiene un índice (que no es de clustering) por el atributo título. Luego, como este atributo no es clave entonces es un índice secundario.

Notemos también que en la consulta tenemos una conjunción (AND) de dos condiciones.

¿Qué pasa si se aplica primero la primer condición?. Como tenemos un atributo que tiene un índice secundario asociado, se aplica primero esta condición y luego se selecciona del resultado a aquellas tuplas que cumplen con la segunda condición sin costo adicional.

Como tenemos un índice, podemos utilizar el método de Index Scan. Luego, el costo de este método considerando una búsqueda con índice secundario se calcula como

$$\text{Costo(Selección)} = \text{Height}(I(A_i, R)) + \lceil \frac{n(R)}{V(A_i, R)} \rceil$$

$$\text{Costo(Selección)} = \text{Height}(I(\text{titulo}, \text{profesores})) + \lceil \frac{n(\text{profesores})}{V(\text{titulo}, \text{profesores})} \rceil$$

$$\text{Costo(Selección)} = 2 + \left\lceil \frac{500}{10} \right\rceil = 2 + 50 = 52$$

¿Qué pasa si se aplica primero la segunda condición?. Como tenemos un atributo simple (no tiene un índice asociado) podemos utilizar el método de File Scan. Luego, el costo de este método se calcula como

$$\text{Costo(Selección)} = B(R) = \frac{n(R)}{F(R)} = \frac{n(\text{profesores})}{F(\text{profesores})} = 20$$

Luego, el costo de aplicar primero la segunda condición es de 20 bloques accedidos.

Si ahora aplicamos la segunda condición tenemos que el costo será $2 + 20 = 22$ bloques accedidos.

Por lo tanto, el método de acceso más eficiente para resolver la consulta es primero aplicar la segunda condición (File Scan) y luego la primera condición (Index Scan) tal que el costo de la consulta es de 22 bloques accedidos.

Item b

Para estimar el tamaño de una selección de la forma $\sigma_{A_i=c}(R)$, utilizaremos la variabilidad de A_i en $R(V(A_i, R))$, que es la cantidad de valores distintos que puede tomar el atributo A_i en dicha relación. Luego, realizaremos la siguiente estimación

$$n(\sigma_{A_i=c}(R)) = \frac{n(R)}{V(A_i, R)}$$

$$n(\text{Selección}) = \frac{n(\text{profesores})}{V(\text{titulo}, \text{profesores}) \times V(\text{genero}, \text{profesores})} = \frac{500}{10 \times 2} = 25$$

Ejercicio 2

Item a

Esquema
$\text{megusta}(\underline{\text{id_usuario}}, \text{id_publicación}, \text{fecha_hora})$

Consulta
$\text{megusta} \bowtie_{\substack{\text{id_usuario} \neq \text{id_usuario}' \wedge \\ \text{id_publicacion} = \text{id_publicacion}'}} \text{megusta}'$

Información
<ul style="list-style-type: none">▪ $n(\text{megusta}) = 100,000,000$ y $F(\text{megusta}) = 1,000$▪ $V(\text{id_usuario}, \text{megusta}) = 50,000$ y $V(\text{id_publicación}, \text{megusta}) = 10,000,000$▪ No se cuenta con índices y se dispone de $M = 1,001$ bloques de memoria disponibles.

Esta consulta tiene como resultado los usuarios (sin incluir al dueño de la publicación) que le dieron me gusta a una publicación.

Notemos que no tenemos índices. Luego, podemos usar el método de loops anidados por bloque, el método de sort-merge y el método de junta hash (variante GRACE).

Notemos también que en la consulta tenemos una conjunción (AND) de dos condiciones.

Comencemos analizando el método de loops anidados por bloque.

¿Cuál es la cantidad de bloques de almacenamiento que ocupa la tabla?.

$$B(R) = \frac{n(R)}{F(R)} = \frac{n(megusta)}{F(megusta)} = \frac{100.000.000}{1.000} = 100.000$$

Notemos que tenemos más bloques en la relación, que bloques de memoria disponibles.

Luego, como en nuestro caso nos entraría CASI la primer tabla entera (salvo por su ultimo bloque), el costo de este método es

$$Costo(R * S) = \min(B(R), B(S)) + \left\lceil \frac{\min(B(R), B(S))}{M - 2} \right\rceil \times B(S)$$

$$Costo(megusta * megusta') = B(megusta) + \left\lceil \frac{B(megusta)}{999} \right\rceil \times B(megusta')$$

$$Costo(megusta * megusta') = B(megusta) + \left\lceil \frac{B(megusta)}{999} \right\rceil \times B(megusta')$$

$$Costo(megusta * megusta') = 100.000 + \left\lceil \frac{100.000}{999} \right\rceil \times 100.000$$

$$Costo(megusta * megusta') = 100.000 + 101 \times 100.000 = 10.200.000$$

Analizamos el método de sort-merge.

Notemos que los archivos no entran en memoria, entonces debe utilizarse un algoritmo de sort externo.

El costo de ordenar R y volverlo a guardar en disco ordenado es

$$2 \times B(R) \times \lceil \log_{M-1}(B(R)) \rceil$$

Una vez ordenados, se hace un merge de ambos archivos que sólo selecciona aquellos pares de tuplas en que coinciden los atributos de junta. El merge recorre una única vez cada archivo, con un costo de $B(R) + B(S)$. El costo total es entonces:

$$\text{Costo}(R * S) = B(R) + B(S) + 2 \times B(R) \times \lceil \log_{M-1}(B(R)) \rceil + 2 \times B(S) \times \lceil \log_{M-1}(B(S)) \rceil$$

Sin embargo, como en nuestro caso tenemos que $R = S$ el costo total es

$$\text{Costo}(R * R') = B(R) + 2 \times B(R) \times \lceil \log_{M-1}(B(R)) \rceil$$

$$\begin{aligned} \text{Costo}(\text{megusta} * \text{megusta}') = \\ B(\text{megusta}) + 2 \times B(\text{megusta}) \times \lceil \log_{1.001-1}(B(\text{megusta})) \rceil \end{aligned}$$

$$\text{Costo}(\text{megusta} * \text{megusta}') = 100.000 + 2 \times 100.000 \times \lceil \log_{1.000}(100.000) \rceil$$

$$\text{Costo}(\text{megusta} * \text{megusta}') = 100.000 + 2 \times 100.000 \times 2 = 500.000$$

Analizamos el método de junta hash (variante GRACE).

Este método cuenta con tres reglas:

- 1) Tenemos que m (cantidad de particiones de hash) debe ser menor o igual a $M-1$ (donde M es la memoria disponible) por la etapa de particionamiento.
- 2) Debe cumplirse que $\frac{B(R)}{m} \leq M - 2$, siendo R la relación más pequeña, por la etapa de loops anidados.
- 3) $m \leq V(A, R)$, $m \leq V(A, S)$ asumiendo A el atributo de junta.

En nuestro caso, se debe cumplir que:

- 1) $m \leq 1000$
- 2) $\frac{100.000}{m} \leq 999$
- 3) $m \leq 50.000$

Luego, el costo de este método es

$$\text{Costo}(R * S) = 3 \times (B(R) + B(S))$$

pero como el join se hace sobre la misma tabla el costo resulta ser

$$\text{Costo}(R * R') = 3 \times B(R)$$

$$\text{Costo}(\text{megusta} * \text{megusta}') = 3 \times B(\text{megusta})$$

$$\text{Costo}(\text{megusta} * \text{megusta}') = 3 \times 100.000 = 300.000$$

Nota: El costo total resulta ser el mencionado ya que tenemos un $B(R)$ para leer bloque a bloque la tabla, un $B(R)$ para grabar las particiones que se hacen una única vez por el atributo de junta y un ultimo $B(R)$ para hacer la parte de loops anidados donde se levanta cada partición entera en memoria y con dos punteros distintos se prueban todas las combinaciones entre filas para ser si cumplen o no la condición de la junta.

Por lo tanto, el método de acceso más eficiente para resolver la consulta es utilizando el método de junta hash (variante GRACE) tal que el costo de la consulta es de 300.000 bloques accedidos.

Item b

Realizaremos la siguiente estimación

$$n(R \bowtie_{R.A=S.A} S) = \frac{n(R) n(S)}{\max(V(A, R), V(A, S))}$$

$$n(megusta * megusta') = \frac{n(megusta) n(megusta')}{\max(V(id_usuario, megusta), V(id_publicacion, megusta))}$$

$$n(megusta * megusta') = \frac{100.000.000 \times 100.000.000}{\max(50.000, 10.000.000)} = \frac{100.000.000 \times 100.000.000}{10.000.000}$$

$$n(megusta * megusta') = 1.000.000.000$$