

CS 525 – QUIZ II
SHANMUGA PRIYA ELLAPPAN
A20465677

Question 1.1.1 Estimate Result Size (5 Points)

Estimate the number of result tuples for the query $q = \sigma_{\text{instructor}=\text{Bob}}(\text{course})$ using the first assumption presented in class (values used in queries are uniformly distributed within the active domain).

$$\begin{aligned} T(q) &= T(\text{course})/V(\text{course}, \text{instructor}) \\ &= 1000/500 \\ &= 2 \text{ tuples} \end{aligned}$$

Question 1.1.2 Estimate Result Size (5 Points)

Estimate the number of result tuples for the query $q = \sigma_{\text{major}=\text{CS} \vee \text{major}=\text{EE}}(\text{students})$ using the first assumption presented in class.

$$\begin{aligned} T(q) &= T(\text{student})/V(\text{student}, \text{major})_{\text{CS}} + T(\text{student})/V(\text{student}, \text{major})_{\text{EE}} \\ &= 40,000/50 + 40,000/50 \\ &= 800 + 800 \\ &= 1600 \text{ Tuples} \end{aligned}$$

Question 1.1.3 Estimate Result Size (7 Points)

Estimate the number of result tuples for the query $q = \sigma_{\text{credits} \geq 29 \wedge \text{credits} \leq 40}(\text{course})$ using the first assumption presented in class.

$$\begin{aligned} T(q) &= [(40 - 29 + 1) * T(\text{course})] / [\max(\text{course}, \text{credits}) - \min(\text{course}, \text{credits}) + 1] \\ &= [12 * 1000] / [60 - 0 + 1] \\ &= 197 \text{ Tuples} \end{aligned}$$

Question 1.1.4 Estimate Result Size (8 Points)

Estimate the number of result tuples for the query q below using the first assumption presented in class.

$q = \sigma_{\text{credits} \leq 30}(\text{student}) \bowtie_{\text{CWID}=\text{student}} \text{registered} \bowtie_{\text{course}=\text{title}} \text{course}$

q1 = $\sigma_{\text{credits} \leq 30}(\text{student})$

$$\begin{aligned} T(q1) &= [(30 - \min(\text{student, credits}) + 1) * T(\text{student})] / [\max(\text{student, credits}) - \min(\text{student, credits}) + 1] \\ &= [(30 - 0 + 1) * 40,000] / [60 - 0 + 1] \\ &= (31 * 40,000) / 61 \\ &= 20,328 \text{ Tuples} \end{aligned}$$

q2 = q1 $\bowtie_{\text{CWID}=\text{student}}$ registered

$$\begin{aligned} T(q2) &= T(q1) * T(\text{registered}) / \max[V(q1, \text{CWID}), V(\text{registered, student})] \\ &= 20,328 * 20,000 / \max(20328, 5000) \\ &= 20,000 \text{ Tuples} \end{aligned}$$

q = q2 $\bowtie_{\text{course}=\text{title}}$ course

$$\begin{aligned} T(q) &= T(q2) * T(\text{course}) / \max[V(q2, \text{course}), V(\text{course, title})] \\ &= 20,000 * 1000 / \max(100, 1000) \\ &= 20,000 \text{ Tuples} \end{aligned}$$

Question 1.2.1 External Sorting (5 Points)

You have a block of size of 128KB and 2GB memory. Compute the maximum size of file that can be sorted by 3-Pass Multiway Sort.

Memory in KB = $2 * 1024 * 1024 = 2,097,152$ KB

$$M = 2,097,152 / 128$$

$$= 16,384 \text{ buffers}$$

$$B(R) \leq M * (M-1)^2 \text{ blocks}$$

$$\leq 16,384 * (16,384 - 1)^2 * 128 \text{ KB}$$

$$\leq 32,766 \text{ GB}$$

Question 1.2.2 I/O Cost Estimation (25 Points)

Consider two relations R and S with $B(R) = 400,000$ and $B(S) = 2,000$ blocks, and $S(R) = 1/40$, and $S(S) = 1/10$. You have $M = 101$ memory pages available. Let R has an index on the joining attribute C. Compute the minimum number of I/O operations needed to join these two relations using tuple-based-nested-loop join (relations are not clustered), block-nested-loop join (relations are clustered), merge-join (the inputs are not sorted but clustered), and non-clustering index-join (read of clustered relation S with uniform distribution assumption and expected 20 matching tuples on R for S) and hash-join (relations are clustered but not sorted). You can assume that the hash function evenly distributes keys across buckets. Justify your result by showing the I/O cost estimation for each join method.

$$T(R) = 400,000 * 40 = 16,000,000$$

$$T(S) = 2000 * 10 = 20,000$$

- Tuple based nested loop join:

S is smaller, thus let's consider it as outer relation

$$\begin{aligned} \text{Total Cost} &= T(S) * (1 + T(R)) \\ &= 20,000 * (1 + 16,000,000) \\ &= 320,000,020,000 \text{ I/Os} \end{aligned}$$

- Block nested loop join:

S is smaller. Thus, let's consider 100 buffers for S and 1 buffer for R

Read chunk of S = 100 I/Os

Read of R = 400,000 I/Os

$$\begin{aligned} \# \text{ of S chunks} &= B(S) / M-1 \\ &= 2000/100 \\ &= 20 \end{aligned}$$

$$\begin{aligned} \text{Total cost} &= 20 * (100 + 400,000) \\ &= 8,002,000 \text{ I/Os} \end{aligned}$$

- Merge Join:

$$\text{Sort cost of R (3 Pass Multiway join algorithm)} = 6 * B(R)$$

$$\text{Sort cost of S (2 Pass Multiway join algorithm)} = 4 * B(S)$$

$$\text{Total cost} = \text{sort cost} + \text{join cost}$$

$$= 7 * B(R) + 5 * B(S)$$

$$= 7 * 400,000 + 5 * 2000$$

$$= 2,800,000 + 10,000$$

$$= 2,810,000 \text{ I/Os}$$

- Non clustering Index join:

$$T(S) + [T(S). T(R) / V(R,C)] = 20,000 + [20,000 * 20]$$

$$= 20,000 + 400,000$$

$$= 420,000 \text{ I/Os}$$

- Hash-Join (relations are clustered but not sorted):

$$3 (B(R) + B(S)) = 3 (400,000 + 2000)$$

$$= 1,206,000 \text{ I/Os}$$