

1. unclustered, height = 3
- ① search down to find the corresponding key: 4 I/O <sup>care.</sup> check if they same with new entry
- ② at most 6 entry in same key: 6 I/O (unclustered)
- ③ Add the new entry: 1 I/O (choose empty, page already in)
- ④ update the node: 1 I/O
- ⇒ 12 I/O
- clustered, height = 5
- ① search: 6 I/O
- ② find the entry: 1 I/O (clustered)
- ③ Add new entry: 1 I/O
- ④ update node: 1 I/O
- ⇒ 9 I/O

## CS150A Homework 2 – Writing

School of Information Science and Technology  
October 29, 2022

October 29, 2022

unclustered: 1 I/O per record  
clustered: 1 I/O per page of record

### 1 B+ tree Refinement (10 pts)

How many I/Os would it cost to insert an entry into our table if we had a height 3, unclustered alternative 3 B+ tree in the **worst case**? Assume that the cache is empty at the beginning and there are at most 6 entries with the same key. And what is the answer for a height 5, clustered alternative 3 B+ tree in the best case? Assume each page is at most 2/3 full.

### 2 Relational Algebra (40 pts)

As shown in the following figures, there are three instances corresponding to three relations: Sailors, Reserves and Boats.

S

sid	sname	rating	age
22	Dustin	7	45.0
29	Brutus	1	33.0
31	Lubber	8	55.0
32	Andy	3	25.0
58	Rusty	10	35.0
74	Horatio	9	35.0

Fig. 1. An instance of Sailors

R

sid	bid	day
22	101	10/10/98
22	104	10/7/98
31	102	11/10/98
31	103	11/6/98
64	102	9/8/98
74	103	9/8/98

Fig. 2. An instance of Reserves

B

bid	bname	color
101	Interlake	blue
102	Interlake	red
103	Clipper	green
104	Marine	red

Fig. 3. An instance of Boats

projection already ensure distinct.

Use relational algebra to describe the following queries.

1. Find the names of sailors who reserved boats after 10/7/98.

$$\pi_{\text{name}}(\sigma_{\text{day} > 10/7/98} (\text{Sailors} \bowtie_{\text{sid}=\text{sid}} \text{Reserves}))$$

2. Find the age of sailors who have reserved boat Interlake with blue color.

$$\pi_{\text{age}}(\text{Sailors} \bowtie_{\text{sid}=\text{sid}} \text{Reserves} \bowtie_{\text{bid}=\text{bid}} (\sigma_{\text{bname} = \text{Interlake}} (\sigma_{\text{color} = \text{blue}} (\text{Boats}))))$$

3. Find the name of boats which have been reserved by sailors rating 7 or higher.

$$\pi_{\text{name}}(\text{Boats} \bowtie_{\text{bid}=\text{bid}} (\text{Reserves} \bowtie_{\text{sid}=\text{sid}} (\sigma_{\text{rating} \geq 7} (\text{Sailors}))))$$

4. Find the color of boats which once reserved by sailors Dustin on 10/10/98.

$$\pi_{\text{color}}(\text{Boats} \bowtie_{\text{bid}=\text{bid}} (\sigma_{\text{day} = 10/10/98} (\text{Reserves} \bowtie_{\text{sid}=\text{sid}} (\sigma_{\text{name} = \text{Dustin}} (\text{Sailors}))))))$$

5. Find the name of sailors who didn't reserved any boat named Marine.

$$\pi_{\text{name}}(\text{Sailors}) - \pi_{\text{name}}(\text{Sailors} \bowtie_{\text{sid}=\text{sid}} (\text{Reserves} \bowtie_{\text{bid}=\text{bid}} (\sigma_{\text{bname} = \text{Marine}} (\text{Boats}))))$$

### 3 Joins (10 pts)

Determine whether each of the following statements is True or False.

- a. Block Nested Loops join will always perform at least as well as Page Nested Loops Join when it comes to minimizing I/Os. **T**

- b. Grace hash join is usually the best algorithm for joins in which the join condition includes an inequality. **F**

require equal

### 4 Query Optimization (40 pts)

Consider two relations R(a, b) and S(a), with 1000 tuples and 500 tuples respectively. We have an index on R.a with 50 unique integer values uniformly distributed in the range [1, 50], and an index on S.a with 25 unique integer values uniformly distributed in the range [1, 25]. We do not have an index on R.b.

index

= clustered

Use selectivity estimation to estimate the number of tuples produced by the following queries.

1. SELECT \* FROM R

$$1000$$

2. SELECT \* FROM R WHERE a = 99

$$1000 \times \frac{1}{50} = 20$$

3. SELECT \* FROM R WHERE b = 99

no index on b.

default selectivity:  $\frac{1}{2}$

$$1000 \times \frac{1}{2} = 500$$

$a: \text{uniform}[1, 50]$        $a \leq 24 \rightarrow a < 25$

4. SELECT \* FROM R WHERE  $a \leq 24$

$$\text{Selectivity} = \frac{24-1}{50-1+1} + \frac{1}{50} = \frac{24}{50} \quad (1000 \times \frac{24}{50} = 480)$$

5. SELECT \* FROM R WHERE  $b \leq 10$

$$\text{Selectivity} = \frac{1}{10} \quad 1000 \times \frac{1}{10} = 100$$

6. SELECT \* FROM R WHERE NOT  $a \leq 24$

where not  $a \leq 24 \rightarrow a \geq 25$

$$\text{Selectivity} = \frac{50-25}{50} + \frac{1}{50} = \frac{26}{50} \quad 1000 \times \frac{26}{50} = 520$$

7. SELECT \* FROM R WHERE  $a \leq 24$  AND  $b \leq 10$

$$\text{Selectivity}(a \leq 24) * \text{Selectivity}(b \leq 10) = \frac{24}{50} \times \frac{1}{10} = \frac{12}{250}$$

$$1000 \times \frac{12}{250} = 48$$

8. SELECT \* FROM R WHERE  $a \leq 24$  OR  $b \leq 10$

$$\text{Selectivity}(a \leq 24) + \text{Selectivity}(b \leq 10) - \text{Selectivity}(a \leq 24) * \text{Selectivity}(b \leq 10) = \frac{24}{50} + \frac{1}{10} - \frac{24}{50} \times \frac{1}{10}$$

$$= \frac{133}{250}$$

$$1000 \times \frac{133}{250} = 532$$

9. SELECT \* FROM R WHERE  $a = b$

$$\text{Selectivity} = \min\left(\frac{1}{50}, \frac{1}{10}\right) = \frac{1}{50} \quad 1000 \times \frac{1}{50} = 20$$

10. SELECT \* FROM R, S WHERE  $R.a = S.a$

$$1000 \times 500 \times \frac{1}{50} = 10000$$

## 11 Appendix - Selectivity Values

- $|column|$  is the number of distinct values in the column

### 11.1 Equalities

Predicate	Selectivity	Assumption
$c = v$	$1 / (\# \text{ of distinct values of } c \text{ in index})$	We know $ c $ .
$c = v$	$1/10$	We don't know $ c $ .
$c_1 = c_2$	$1 / \max(\# \text{ of distinct values of } c_1, \# \text{ of distinct values of } c_2)$	We know $ c_1 $ and $ c_2 $ .
$c_1 = c_2$	$1 / (\# \text{ of distinct values of } c_i)$	We know $ c_i $ but not $ \text{other column} $ .
$c_1 = c_2$	$1/10$	We don't know $ c_1 $ or $ c_2 $ .

### 11.2 Inequalities on Integers

Predicate	Selectivity	Assumption
$c < v$	$(v - \min(c)) / (\max(c) - \min(c) + 1)$	We know $\max(c)$ and $\min(c)$ .
$c > v$	$(\max(c) - v) / (\max(c) - \min(c) + 1)$	$c$ is an integer.
$c < v$ $c > v$	$1/10$	We don't know $\max(c)$ and $\min(c)$ . $c$ is an integer.
$c \leq v$	$(v - \min(c)) / (\max(c) - \min(c) + 1) + (1/ c )$	We know $\max(c)$ and $\min(c)$ .
$c \geq v$	$(\max(c) - v) / (\max(c) - \min(c) + 1) + (1/ c )$	$c$ is an integer.
$c \leq v$ $c \geq v$	$1/10$	We don't know $\max(c)$ and $\min(c)$ . $c$ is an integer.

### 11.3 Inequalities on Floats

Predicate	Selectivity	Assumption
$c \geq v$	$(\max(c) - v) / (\max(c) - \min(c))$	We know $\max(c)$ and $\min(c)$ . $c$ is a float.
$c \geq v$	$1/10$	We don't know $\max(c)$ and $\min(c)$ . $c$ is a float.
$c \leq v$	$(v - \min(c)) / (\max(c) - \min(c))$	We know $\max(c)$ and $\min(c)$ . $c$ is a float.
$c \leq v$	$1/10$	We don't know $\max(c)$ and $\min(c)$ . $c$ is a float.