

OPTIMIZATION OF RELATIONAL ALGEBRA

--> Do Selection first, then project, only join when you have to!

- ① JOIN v. CART. PRODUCT
- joins are better

$$R \bowtie S = \sigma_c(R \times S)$$

- ② PUSH SELECTION DOWN
⇒ DO SELECTION AS EARLY AS POSSIBLE

- it reduces the size of the data

$$\sigma_c(\pi_L(R)) \leftrightarrow \pi_L(\sigma_c(R))$$

$$\sigma_{R \times S}(R \bowtie_{A=a.s.b} S) \leftrightarrow \sigma_{R \times S}(R) \bowtie_{A=a.s.b} S$$

- ③ AVOID UNNECESSARY JOINS

find customers having account balances below 100 and loans above 10000

$$R1 \leftarrow \pi_{\text{CustName}}(\text{Depositor} \bowtie \pi_{\text{AccountNum}}(\sigma_{\text{balance} < 100}(\text{Account})))$$

$$R2 \leftarrow \pi_{\text{CustName}}(\text{Borrower} \bowtie \pi_{\text{LoanNum}}(\sigma_{\text{Amount} > 10000}(\text{Loan})))$$

$$\text{Result} \leftarrow R1 \cap R2$$

⇒ Much better than trying to join 4 tables!

RELATIONSHIPS AMONG OPERATORS

- ① JOIN ↔ CARTESIAN PRODUCT + SELECT
 $R \bowtie S = \sigma_c(R \times S)$
- ② SELECTION is commutative
 $\sigma_{c_2}(\sigma_{c_1}(R)) = \sigma_{c_1}(\sigma_{c_2}(R)) = \sigma_{c_1 \text{ AND } c_2}(R)$
- ③ ORDER BETWEEN SELECTION & PROJECTION
 $\sigma_c(\pi_L(R)) \rightarrow \pi_L(\sigma_c(R))$
 $\pi_L(\sigma_c(R)) \rightarrow \sigma_c(\pi_L(R))$ * ONLY if L has all cols needs for C
- ④ JOIN is COMMUTATIVE
 $R \bowtie S = S \bowtie R$
- ⑤ ORDER BETWEEN SELECTION & JOIN
 $\sigma_{c_1}(R \bowtie_{c_2} S) \rightarrow (\sigma_{c_1}(R) \bowtie_{c_2} S)$

- ① Set operations: UNION, INTERSECTION, DIFFERENCE



- ② PROJECTION $\pi_L(R)$: SELECT clause - chooses some columns to use
L is a list consisting of col name, col renames (eg A as B), or expressions (eg A+B as Z)

- ③ SELECTION $\sigma_c(R)$: WHERE clause - choose some tuples to use
C is a conditional to apply to each tuple in R

- ③ COMBINING TABLES: FROM clause - choose which tables to use & how to join

- ① CROSS-PRODUCT $A \times B$: pairs each $a \in A$ with each $b \in B$ ⇒ makes huge tables
- ② NATURAL JOIN $A \bowtie B$: pairs each matching attribute in matching columns
- ③ THETA JOIN $A \bowtie_{c} B$: pairs each $a \in A$ with each $b \in B$ if $c(a,b)$ holds true
- ④ OUTER JOIN $A \bowtie_{\perp} B$: pairs according to theta join, then pass the dangling tuples with \perp

A:	X	Y
	0	1
	1	0

B:	X	C
	1	0
	1	1

$A \times B$:	X	Y	B	X	C
	0	1	1	0	0
	1	0	1	0	0
	0	1	1	1	1
	1	0	1	1	1

$A \bowtie B$:	X	Y	C
	1	0	0
	1	0	1

$A \bowtie_{Y=C} B$:	X	Y	C
	0	1	1
	1	0	1

$A \bowtie_{\perp} B$:	X	Y	C
	1	0	0
	1	0	1
	0	1	\perp

- ④ RENAMING: $\rho_{s(A_1, A_2, \dots, A_n)}(R)$: AS operator

- ⑤ DUPLICATE ELIMINATION: $\delta(R)$: DISTINCT operator - returns R with one copy of each tuple in R
⇒ Turns a bag into a set

- ⑥ SORTING: $\tau_L(R)$: ORDER BY clause

- L is a list of fields to sort by, where ties are broken by fields later in the list

- ⑦ AGGREGATION & GROUPING: $\gamma_L(R)$: GROUP BY clause

- L is a list of GROUPING ATTRIBUTES (attributes $\in R$) and AGGREGATED ATTRIBUTES (operator applied to cols of R)

$\gamma_{\text{activity, sum(hours)} \rightarrow \text{time spent}}(\text{HoursLog})$

- executing $\gamma_L(R)$:
- ① Partition R into groups, where each group has tuples with a distinct assignment of the grouping attributes
→ if no grouping attrs specified, R is one group
 - ② For each group, produce one tuple consisting of:
 - the grouping attributes for that group
 - the aggregations over the tuples in that group