

Exercise

Let's consider the following relational schema for a group of insurance companies located in different cities:

CUSTOMER(Id_cust, Name, Age, City_cust)

INSURANCE_COMPANY(Id_company, Id_Director, nEmployee, City)

POLICY(Id_policy, Id_cust, Id_company, expiry_date)

Primary keys are underlined in the relations. Moreover, Id_cust in POLICY is foreign key of CUSTOMER; Id_company in POLICY is foreign key of INSURANCE_COMPANY and Id_Director in INSURANCE_COMPANY foreign key of CUSTOMER.

A customer can have more than one policy in the same company or in different companies.

Expiry_date in POLICY is a year.

Assume that:

$n_{\text{CUSTOMER}} = 2000$

$n_{\text{INSURANCE_COMPANY}} = 20$

$n_{\text{POLICY}} = 100.000$

$V(\text{Id_cust}, \text{POLICY}) = 2000$

$V(\text{Id_company}, \text{POLICY}) = 20$

$V(\text{expiry_date}, \text{POLICY}) = 20$

$V(\text{City}, \text{INSURANCE_COMPANY}) = 5$

Given the query:

“Name of customers holding policies with companies located in Pisa and with expiry date 2010”

- 1) express the query as a relational-algebra expression;
- 2) show the basic steps of the query optimization process in terms of relational-algebra expression transformations
- 3) give an efficient strategy for computing the query.

Point 1

$\Pi_{\text{CUSTOMER.Name}} (\sigma_{\text{INSURANCE_COMPANY.City=Pisa and POLICY.expiry_date=2010}} ($
 $(\text{CUSTOMER} \bowtie_{\text{CUSTOMER.Id_cust=POLICY.Id_cust}} \text{POLICY})$
 $\bowtie_{\text{POLICY.Id_company=INSURANCE_COMPANY.Id_company}} \text{INSURANCE_COMPANY}))$

Let C, IC and P denote CUSTOMER, INSURANCE_COMPANY and POLICY, respectively.

$\Pi_{\text{C.Name}} (\sigma_{\text{IC.City=Pisa and P.expiry_date=2010}} ((\text{C} \bowtie_{\text{C.Id_cust=P.Id_cust}} \text{P}) \bowtie_{\text{P.Id_company=IC.Id_company}} \text{IC}))$

Point 2

$\sigma_{\text{IC.City=Pisa and P.expiry_date=2010}} (.....)$ can be rewritten as: $\sigma_{\text{IC.City=Pisa}} (\sigma_{\text{P.expiry_date=2010}} (.....))$

$\Pi_{\text{C.Name}} (\sigma_{\text{IC.City=Pisa}} (\sigma_{\text{P.expiry_date=2010}} ((\text{C} \bowtie_{\text{C.Id_cust=P.Id_cust}} \text{P}) \bowtie_{\text{P.Id_company=IC.Id_company}} \text{IC})))$

Push selection down

$\Pi_{\text{C.Name}} ((\text{C} \bowtie_{\text{C.Id_cust=P.Id_cust}} (\sigma_{\text{P.expiry_date=2010}} (\text{P}))) \bowtie_{\text{P.Id_company=IC.Id_company}} (\sigma_{\text{IC.City=Pisa}} (\text{IC})))$

Push projection down

$\Pi_{\text{C.Name}} ((\Pi_{\text{C.Name, C.Id_cust}} \text{C}) \bowtie_{\text{C.Id_cust=P.Id_cust}} (\Pi_{\text{P.Id_cust, P.Id_company}} (\sigma_{\text{P.expiry_date=2010}} \text{P}))$
 $\bowtie_{\text{P.Id_company=IC.Id_company}} (\Pi_{\text{IC.Id_company}} (\sigma_{\text{IC.City=Pisa}} (\text{IC})))$

We evaluate the size and the number of different values for the new relations.

Let $C' = \Pi_{\text{C.Name, C.Id_cust}} (\text{C})$

$n_{C'} = n_{\text{CUSTOMER}} = 2000$ Id_cust is a key

Let $P' = \sigma_{P.\text{expiry_date}=2010}(P)$

$$n_{P'} = n_{\text{POLICY}} / V(\text{expiry_date}, \text{POLICY}) = (100.000/20) = 5.000$$

$$V(\text{Id_cust}, P') = \min(n_{P'}, V(\text{Id_cust}, P)) = \min(5.000, 2.000) = 2.000$$

$$V(\text{Id_company}, P') = \min(n_{P'}, V(\text{Id_company}, P)) = \min(5.000, 20) = 20$$

Let $P'' = \Pi_{P.\text{Id_cust}, P.\text{Id_company}}(P')$

$$n_{P''} = \min(n_{P'}, V(\text{Id_cust}, P') * V(\text{Id_company}, P')) = \min(5.000, 2.000 * 20) = 5.000$$

$$V(\text{Id_cust}, P'') = 2.000$$

$$V(\text{Id_company}, P'') = 20$$

Let $IC' = \sigma_{IC.\text{City}=Pisa}(IC)$

$$n_{IC'} = (n_{\text{INSURANCE_COMPANY}} / V(\text{City}, \text{INSURANCE_COMPANY})) = (20/5) = 4$$

$$V(\text{Id_company}, IC') = n_{IC'} = 4$$

Let $IC'' = \Pi_{IC.\text{Id_company}}(IC')$

$$n_{IC''} = n_{IC'} = 4 \quad (\text{Id_company is a key})$$

Point 3

The query expression can be rewritten using natural join operator. Natural join is commutative.

$$\Pi_{C.\text{Name}}(C' \bowtie P'' \bowtie IC'')$$

We estimate the size of different combinations of join.

Let $T1 = (C' \bowtie_{C'.\text{Id_cust}=P''.\text{Id_cust}} P'')$

Number of records in the result:

Id_cust in P'' is foreign key of C' (note that C' and C have the same values of Id_cust)

$$n_{T1} = n_{P''} = 5000$$

Let $T2 = (C' \bowtie IC'')$ Cartesian product

Number of records in the result:

$$n_{T2} = (n_{C'} * n_{IC''}) = 2000 * 4 = 8000$$

Let $T3 = (P'' \bowtie_{P''.\text{Id_company}=IC''.\text{Id_company}} IC'')$

Number of records in the result:

Id_company in P'' is not foreign key of IC''

Id_company in P'' is a key of IC''

$$n_{T3} < n_{P''} < 5.000$$

More precisely :

$n_{T3} = \text{number of policies at each insurance company} * \text{number of companies}$

-number of policies at each insurance company: $n_{P''} / V(\text{Id_company}, P'') = 5.000 / 20 = 250$

-number of insurance companies: $n_{IC''} = 4$

$$n_{T3} = 250 * 4 = 1.000$$

Rule applied by the optimizer:

$$\min(n_{P''} * (n_{IC''} / V(\text{Id_company}, IC'')), n_{IC''} * (n_{P''} / V(\text{Id_company}, P''))) =$$

$$\min(5000 * (4/4), 4 * (5.000 / 20)) = \min(5.000, 1.000) = 1.000$$

The best ordering of join is : $(C' \bowtie_{C'.\text{Id_cust}=P''.\text{Id_cust}} (P'' \bowtie_{P''.\text{Id_company}=IC''.\text{Id_company}} IC''))$

An efficient strategy for solving the query is:

$$\Pi_{C.\text{Name}}(C' \bowtie_{C'.\text{Id_cust}=P''.\text{Id_cust}} (P'' \bowtie_{P''.\text{Id_company}=IC''.\text{Id_company}} IC''))$$