Extended Relational-Algebra-Operations

- Generalized Projection
- Aggregate Functions
- Outer Join

Generalized Projection

Extends the projection operation by allowing arithmetic functions to be used in the projection list.

 $\prod_{F_1},_{F_2},\ldots,_{F_n}(E)$

- *E* is any relational-algebra expression
- Each of F_1 , F_2 , ..., F_n are are arithmetic expressions involving constants and attributes in the schema of *E*.
- Given relation credit_info(customer_name, limit, credit_balance), find how much more each person can spend:

∏*customer_name, limit – credit_balance* (credit_info)

Aggregate Functions and Operations

Aggregation function takes a collection of values and returns a single value as a result.

avg: average value
min: minimum value
max: maximum value
sum: sum of values
count: number of values

Aggregate operation in relational algebra

$$_{G_1,G_2,\ldots,G_n} \mathcal{G}_{F_1(A_1),F_2(A_2),\ldots,F_n(A_n)}(E)$$

E is any relational-algebra expression

- G_1, G_2, \ldots, G_n is a list of attributes on which to group (can be empty)
- Each F_i is an aggregate function
- Each A_i is an attribute name
- A_1, \ldots, A_n are disjoint from G_1, \ldots, G_n

Aggregate Operation – Example

Relation *r*.

A	В	С
α	α	7
α	β	7
β	β	3
β	β	10



Aggregate Operation – Example

Relation *account* grouped by *branch-name*:

branch_name	account_number	balance
Perryridge	A-102	400
Perryridge	A-201	900
Brighton	A-217	750
Brighton	A-215	750
Redwood	A-222	700

branch_name $g_{sum(balance)}(account)$

branch_name	sum(balance)
Perryridge	1300
Brighton	1500
Redwood	700

Aggregate Functions (Cont.)

Result of aggregation does not have a name

- Can use rename operation to give it a name
- For convenience, we permit renaming as part of aggregate operation

branch_name *g* sum(balance) as sum_balance (account)

Outer Join

- An extension of the join operation that avoids loss of information.
- Computes the join and then adds tuples from one relation that does not match tuples in the other relation to the result of the join.
- Uses *null* values:
 - *null* signifies that the value is unknown or does not exist
 - All comparisons involving *null* are (roughly speaking) false by definition.
 - > We shall study precise meaning of comparisons with nulls later

Banking Example

branch (branch_name, branch_city, assets)

customer (customer_name, customer_street, customer_city)

account (account_number, branch_name, balance)

loan (loan_number, branch_name, amount)

depositor (customer_name, account_number)

borrower (customer_name, loan_number)

Outer Join – Example

Relation *loan*

loan_number	branch_name	amount
L-170	Downtown	3000
L-230	Redwood	4000
L-260	Perryridge	1700

Relation *borrower*

customer_name	loan_number
Jones	L-170
Smith	L-230
Hayes	L-155

Outer Join – Example

Join

 $loan \Join borrower$

loan_number	branch_name	amount	customer_name
L-170	Downtown	3000	Jones
L-230	Redwood	4000	Smith

Left Outer Join

loan ____ borrower

loan_number	branch_name	amount	customer_name
L-170 L-230 L-260	Downtown Redwood Perryridge	3000 4000 1700	Jones Smith

Outer Join – Example

Right Outer Join

loan \bowtie borrower

loan_number	branch_name	amount	customer_name
L-170	Downtown	3000	Jones
L-230 L-155	null	4000 null	Hayes

Full Outer Join

loan⊐⋈⊂ *borrower*

loan_number	branch_name	amount	customer_name
L-170	Downtown	3000	Jones
L-230	Redwood	4000	Smith
L-260	Perryridge	1700	null
L-155	null	null	Hayes

Null Values

- It is possible for tuples to have a null value, denoted by *null*, for some of their attributes
- *null* signifies an unknown value or that a value does not exist.
- The result of any arithmetic expression involving *null* is *null*.
- Aggregate functions simply ignore null values (as in SQL)
- For duplicate elimination and grouping, null is treated like any other value, and two nulls are assumed to be the same (as in SQL)

Null Values

Comparisons with null values return the special truth value: unknown

- If *false* was used instead of *unknown*, then *not* (A < 5) would not be equivalent to $A \ge 5$
- *Three-valued* logic using the truth value *unknown*:
 - OR: (unknown or true) = true, (unknown or false) = unknown (unknown or unknown) = unknown
 - AND: (true and unknown) = unknown, (false and unknown) = false, (unknown and unknown) = unknown
 - NOT: (not unknown) = unknown
 - In SQL "P is unknown" evaluates to true if predicate P evaluates to unknown
- Result of select predicate is treated as *false* if it evaluates to *unknown*

Modification of the Database

- The content of the database may be modified using the following operations:
 - Deletion
 - Insertion
 - Update
- All these operations are expressed using the assignment operator.

Deletion

- A delete request is expressed similarly to a query, except instead of displaying tuples to the user, the selected tuples are removed from the database.
- Can delete only whole tuples; cannot delete values on only particular attributes
 - A deletion is expressed in relational algebra by:

 $r \leftarrow r - E$

where r is a relation and E is a relational algebra query.

Banking Example

branch (branch_name, branch_city, assets)

customer (customer_name, customer_street, customer_city)

account (account_number, branch_name, balance)

loan (loan_number, branch_name, amount)

depositor (customer_name, account_number)

borrower (customer_name, loan_number)

Deletion Examples

Delete all account records in the Perryridge branch. $account \leftarrow account - \sigma_{branch_name} = "Perryridge" (account)$

Delete all loan records with amount in the range of 0 to 50 loan \leftarrow loan $-\sigma$ amount \geq 0 and amount \leq 50 (loan)

Delete all accounts at branches located in Needham.

$$r_1 \leftarrow \sigma_{branch_city} = "Needham" (account \Join branch)$$

 $r_2 \leftarrow \Pi_{account_number, branch_name, balance} (r_1)$
 $r_3 \leftarrow \Pi_{customer_name, account_number} (r_2 \Join depositor)$
 $account \leftarrow account - r_2$
 $depositor \leftarrow depositor - r_3$

Insertion

To insert data into a relation, we either:

- specify a tuple to be inserted
- write a query whose result is a set of tuples to be inserted

in relational algebra, an insertion is expressed by:

$$r \leftarrow r \cup E$$

where r is a relation and E is a relational algebra expression.

The insertion of a single tuple is expressed by letting E be a constant relation containing one tuple.

Insertion Examples

Insert information in the database specifying that Smith has \$1200 in account A-973 at the Perryridge branch.

```
account \leftarrow account \cup {("A-973", "Perryridge", 1200)}
depositor \leftarrow depositor \cup {("Smith", "A-973")}
```

Provide as a gift for all loan customers in the Perryridge branch, a \$200 savings account. Let the loan number serve as the account number for the new savings account.

 $r_{1} \leftarrow (\sigma_{branch_name = "Perryridge"}(borrower \bowtie loan))$ account $\leftarrow account \cup \prod_{loan_number, branch_name, 200} (r_{1})$ depositor $\leftarrow depositor \cup \prod_{customer_name, loan_number} (r_{1})$

Updating

- A mechanism to change a value in a tuple without changing all values in the tuple
- Use the generalized projection operator to do this task

$$r \leftarrow \prod_{F_1, F_2, \dots, F_l} (r)$$

- Each *F_i* is either
 - The *i*th attribute of *r*, if the *i*th attribute is not updated, or,
 - if the attribute is to be updated F_i is an expression, involving only constants and the attributes of r, which gives the new value for the attribute

Update Examples

Make interest payments by increasing all balances by 5 percent.

 $account \leftarrow \prod_{account_number, branch_name, balance * 1.05} (account)$

Pay all accounts with balances over \$10,000 6 percent interest and pay all others 5 percent

account $\leftarrow \prod_{account_number, branch_name, balance * 1.06} (\sigma_{BAL > 10000} (account))$ $\cup \prod_{account_number, branch_name, balance * 1.05} (\sigma_{BAL \le 10000} (account))$ (account))

Introduction to SQL (I)

Roadmap of This Lecture

- Overview of The SQL Query Language
- Data Definition
- Basic Query Structure
- Additional Basic Operations
- Set Operations
- Null Values
- Aggregate Functions

History

- IBM Sequel language developed as part of System R project at the IBM San Jose Research Laboratory
- Renamed Structured Query Language (SQL)
- ANSI and ISO standard SQL:
 - SQL-86
 - SQL-89
 - SQL-92
 - SQL:1999 (language name became Y2K compliant!)
 - SQL:2003
- Commercial systems offer most, if not all, SQL-92 features, plus varying feature sets from later standards and special proprietary features.
 - Not all examples here may work on your particular system.

Data Definition Language

The SQL data-definition language (DDL) allows the specification of information about relations, including:

- The schema for each relation.
- The domain of values associated with each attribute.
- Integrity constraints
- And as we will see later, also other information such as
 - The set of indices to be maintained for each relations.
 - Security and authorization information for each relation.
 - The physical storage structure of each relation on disk.

Domain Types in SQL

- **char(n).** Fixed length character string, with user-specified length *n*.
- varchar(n). Variable length character strings, with user-specified maximum length n.
- **int.** Integer (a finite subset of the integers that is machine-dependent).
- smallint. Small integer (a machine-dependent subset of the integer domain type).
- numeric(p,d). Fixed point number, with user-specified precision of p digits, with d digits to the right of decimal point.
- real, double precision. Floating point and double-precision floating point numbers, with machine-dependent precision.
- float(n). Floating point number, with user-specified precision of at least n digits.
- More are covered in "Data types and schemas" later.

Create Table Construct

An SQL relation is defined using the **create table** command:

...,

```
create table r (A_1 D_1, A_2 D_2, ..., A_n D_n, (integrity-constraint_1),
```

(integrity-constraint_k))

- *r* is the name of the relation
- each A_i is an attribute name in the schema of relation r
- D_i is the data type of values in the domain of attribute A_i

Example:

create table instructor (*ID* char(5), *name* varchar(20), *dept_name* varchar(20), *salary* numeric(8,2))

Integrity Constraints in Create Table

not null

- **primary key** $(A_1, ..., A_n)$
- **foreign key** $(A_m, ..., A_n)$ references *r*

Example: Declare *branch_name* as the primary key for *branch*

create table instructor (ID char(5), name varchar(20) not null, dept_name varchar(20), salary numeric(8,2), primary key (ID), foreign key (dept_name) references department);

primary key declaration on an attribute automatically ensures not null

And a Few More Relation Definitions

create table student (

IDvarchar(5),namevarchar(20) not null,dept_namevarchar(20),tot_crednumeric(3,0),primary key (ID),foreign key (dept_name) references department);

create table takes (

ID	varchar(5),
course_id	varchar(8),
sec_id	varchar(8),
semester	varchar(6),
year	numeric(4,0),
grade	varchar(2),
primary key	(ID, course_id, sec_id, semester, year) ,
foreign key (ID) references student,
foreign key	(course_id, sec_id, semester, year) references section);

 Note: sec_id can be dropped from primary key above, to ensure a student cannot be registered for two sections of the same course in the same semester

And more still

create table course (

course_idvarchar(8),titlevarchar(50),dept_namevarchar(20),creditsnumeric(2,0),primary key(course_id),foreign key(dept_name) references department);

Updates to tables

Insert

insert into instructor values ('10211', 'Smith', 'Biology', 66000);

Delete

- delete from student
- Drop Table
 - drop table r
- Alter
 - alter table *r* add *A D*
 - where A is the name of the attribute to be added to relation
 r and D is the domain of A.
 - All tuples in the relation are assigned *null* as the value for the new attribute.
 - alter table *r* drop *A*
 - where A is the name of an attribute of relation r
 - Dropping of attributes not supported by many databases.

Basic Query Structure

A typical SQL query has the form:

select $A_1, A_2, ..., A_n$ **from** $r_1, r_2, ..., r_m$ **where** P

- *A_i* represents an attribute
- r_i represents a relation
- *P* is a predicate.
- The result of an SQL query is a relation.

The select Clause

The **select** clause list the attributes desired in the result of a query

- corresponds to the projection operation of the relational algebra
- Example: find the names of all instructors:

select name from instructor

- NOTE: SQL names are case insensitive (i.e., you may use upper- or lower-case letters.)
 - E.g., $Name \equiv NAME \equiv name$
 - Some people use upper case wherever we use bold font.

The select Clause (Cont.)

- SQL *allows* duplicates in relations as well as in query results.
- To force the elimination of duplicates, insert the keyword distinct after select.
- Find the department names of all instructors, and remove duplicates

select distinct dept_name
from instructor

The keyword **all** specifies that duplicates not be removed.

select all dept_name from instructor

The select Clause (Cont.)

An asterisk in the select clause denotes "all attributes"

select *
from instructor

The select clause can contain arithmetic expressions involving the operation, +, -, *, and /, and operating on constants or attributes of tuples. (This corresponds to generalized projection.)

The query:

select ID, name, salary/12 from instructor

would return a relation that is the same as the *instructor* relation, except that the value of the attribute *salary* is divided by 12.

The where Clause

The where clause specifies conditions that the result must satisfy

- Corresponds to the selection predicate of the relational algebra.
- To find all instructors in Comp. Sci. dept with salary > 80000

select name
from instructor
where dept_name = 'Comp. Sci.' and salary > 80000

- Comparison results can be combined using the logical connectives and, or, and not.
- Comparisons can be applied to results of arithmetic expressions.

The from Clause

The **from** clause lists the relations involved in the query

- Corresponds to the Cartesian product operation of the relational algebra.
- Find the Cartesian product *instructor X teaches*

select *
from instructor, teaches

- generates every possible instructor teaches pair, with all attributes from both relations.
- Cartesian product not very useful directly, but useful combined with where-clause condition (selection operation in relational algebra).

Cartesian Product

instructor

teaches

ID	name	dept n	ame sala	ıru	8	ID	course_i	d sec_	_id sem	ester	year
10101 12121 15151 22222 32343	Srinivasa Wu Mozart Einstein El Said	in Comp Finance Music Physic Histor	. Sci. 650 ce 900 400 cs 950 ry 600)000)000)000)000)000		10101 10101 10101 12121 15151 22222	CS-101 CS-315 CS-347 FIN-201 MU-199 PHY-10	1 1 1 1 1	Fall Spri Fall Spri Spri Fall	ing ing ing	2009 2010 2009 2010 2010 2010 2009
	Inst.ID	name	dept_name	salary	tea	ches.ID	course_id	sec_id	semester	year	
	10101 10101 10101 10101 10101 10101 	Srinivasan Srinivasan Srinivasan Srinivasan Srinivasan Srinivasan 	Comp. Sci. Comp. Sci. Comp. Sci. Comp. Sci. Comp. Sci. Comp. Sci. 	65000 65000 65000 65000 65000 65000 	10 10 10 12 15 22	0101 0101 0101 2121 5151 2222 	CS-101 CS-315 CS-347 FIN-201 MU-199 PHY-101 	1 1 1 1 1 1 	Fall Spring Fall Spring Spring Fall 	2009 2010 2009 2010 2010 2009 	
	12121 12121 12121 12121 12121 12121 12121 	Wu Wu Wu Wu Wu 	Finance Finance Pinance Finance Finance Pinance 	90000 90000 90000 90000 90000 90000 	10 10 10 12 15 22	0101 0101 0101 012121 0151 02222 	CS-101 CS-315 CS-347 FIN-201 MU-199 PHY-101 	1 1 1 1 1 	Fall Spring Fall Spring Spring Fall 	2009 2010 2009 2010 2010 2009 	
	• • •	•••		•••		•••	•••	• • •			38

Joins

For all instructors who have taught courses, find their names and the course ID of the courses they taught.

```
select name, course_id
from instructor, teaches
where instructor.ID = teaches.ID
```

Find the course ID, semester, year and title of each course offered by the Comp. Sci. department

Natural Join

Natural join matches tuples with the same values for all common attributes, and retains only one copy of each common column

select *

from instructor natural join teaches;

ID	name	dept_name	salary	course_id	sec_id	semester	year
10101	Srinivasan	Comp. Sci.	65000	CS-101	1	Fall	2009
10101	Srinivasan	Comp. Sci.	65000	CS-315	1	Spring	2010
10101	Srinivasan	Comp. Sci.	65000	CS-347	1	Fall	2009
12121	Wu	Finance	90000	FIN-201	1	Spring	2010
15151	Mozart	Music	40000	MU-199	1	Spring	2010
22222	Einstein	Physics	95000	PHY-101	1	Fall	2009
32343	El Said	History	60000	HIS-351	1	Spring	2010
45565	Katz	Comp. Sci.	75000	CS-101	1	Spring	2010
45565	Katz	Comp. Sci.	75000	CS-319	1	Spring	2010
76766	Crick	Biology	72000	BIO-101	1	Summer	2009
76766	Crick	Biology	72000	BIO-301	1	Summer	2010
83821	Brandt	Comp. Sci.	92000	CS-190	1	Spring	2009
83821	Brandt	Comp. Sci.	92000	CS-190	2	Spring	2009
83821	Brandt	Comp. Sci.	92000	CS-319	2	Spring	2010
98345	Kim	Elec. Eng.	80000	EE-181	1	Spring	2009

Natural Join Example

List the names of instructors along with the course ID of the courses that they taught.

select name, course_id
 from instructor, teaches
 where instructor.ID = teaches.ID;

select name, course_id
 from instructor natural join teaches;

The Rename Operation

The SQL allows renaming relations and attributes using the **as** clause: old-name **as** new-name

E.g.,

 select ID, name, salary/12 as monthly_salary from instructor

Find the names of all instructors who have a higher salary than some instructor in 'Comp. Sci'.

select distinct T. name
 from instructor as T, instructor as S
 where T.salary > S.salary and S.dept_name = 'Comp. Sci.'

Keyword **as** is optional and may be omitted instructor **as** $T \equiv$ instructor T

String Operations

- SQL includes a string-matching operator for comparisons on character strings. The operator like uses patterns that are described using two special characters:
 - percent (%). The % character matches any substring.
 - underscore (_). The _ character matches any character.
- Find the names of all instructors whose name includes the substring "dar".

select name from instructor where name like '%dar%'

Match the string "100%"

like '100 \%' **escape '**\'

in that above we use backslash (\) as the escape character.

String Operations (Cont.)

- Patterns are case sensitive.
- Pattern matching examples:
 - 'Intro%' matches any string beginning with "Intro".
 - '%Comp' matches any string containing "Comp" as a suffix.
 - '___' matches any string of exactly three characters.
 - '___%' matches any string of at least three characters.
- SQL supports a variety of string operations such as
 - concatenation (using "II")
 - converting from upper to lower case (and vice versa)
 - > upper() and lower()
 - finding string length, extracting substrings, etc.

Ordering the Display of Tuples

List in alphabetic order the names of all instructors

select distinct name from instructor order by name

- We may specify desc for descending order or asc for ascending order, for each attribute; ascending order is the default.
 - Example: order by name desc
- Can sort on multiple attributes
 - Example: order by dept_name, name

Where Clause Predicates

- SQL includes a **between** comparison operator
- Example: Find the names of all instructors with salary between \$90,000 and \$100,000 (that is, \geq \$90,000 and \leq \$100,000)
 - select name
 from instructor
 where salary between 90000 and 100000
- Tuple comparison
 - select name, course_id
 from instructor, teaches
 where (instructor.ID, dept_name) = (teaches.ID, 'Biology');

Duplicates

- In relations with duplicates, SQL can define how many copies of tuples appear in the result.
- Multiset versions of some of the relational algebra operators given multiset relations r_1 and r_2 :
 - 1. $\sigma_{\theta}(r_1)$: If there are c_1 copies of tuple t_1 in r_1 , and t_1 satisfies selections σ_{θ} , then there are c_1 copies of t_1 in $\sigma_{\theta}(r_1)$.
 - 2. $\Pi_{A}(r)$: For each copy of tuple t_{1} in r_{1} , there is a copy of tuple $\Pi_{A}(t_{1})$ in $\Pi_{A}(r_{1})$ where $\Pi_{A}(t_{1})$ denotes the projection of the single tuple t_{1} .
 - 3. $r_1 \ge r_2$: If there are c_1 copies of tuple t_1 in r_1 and c_2 copies of tuple t_2 in r_2 , there are $c_1 \ge c_2$ copies of the tuple t_1 . t_2 in $r_1 \ge r_2$

Duplicates (Cont.)

Example: Suppose multiset relations r_1 (*A*, *B*) and r_2 (*C*) are as follows:

 $r_1 = \{(1, a) (2, a)\}$ $r_2 = \{(2), (3), (3)\}$

Then $\Pi_B(r_1)$ would be {(a), (a)}, while $\Pi_B(r_1) \ge r_2$ would be

 $\{(a,2), (a,2), (a,3), (a,3), (a,3), (a,3)\}$

SQL duplicate semantics:

select $A_{1,,} A_{2}, ..., A_{n}$ from $r_{1}, r_{2}, ..., r_{m}$ where P

is equivalent to the *multiset* version of the expression:

$$\prod_{A_1,A_2,\ldots,A_n} (\sigma_P(r_1 \times r_2 \times \ldots \times r_m))$$

End