

.

SQL

Part II: Advanced Queries

Aggregation

- Significant extension of relational algebra

<u>sid</u>	sname	rating	age
28	yuppy	7	15
31	debby	7	10
22	conny	5	10
58	lilly	10	13

- “Count the number of tuples in Skaters”

```
SELECT COUNT(*)  
FROM skaters
```

count
4

- “Count how many different ratings?”

```
SELECT COUNT(DISTINCT rating)  
FROM skaters
```

count
3

Result is a relation with only one tuple



Aggregation

- Syntax: **COUNT**, **SUM**, **AVG**, **MAX**, **MIN** apply to single attribute/column.
- Additionally, **COUNT (*)**
- “What is the average age of skaters with rating 7?”

```
SELECT AVG(age)
FROM skaters
WHERE rating = 7
```
- What is the average age of skaters with rating 7, and how many are there?”

```
SELECT AVG(age) , COUNT(*)
FROM skaters
WHERE rating = 7
```

3

<u>sid</u>	sname	rating	age
28	yuppy	7	15
31	debby	7	10
22	conny	5	10
58	lilly	10	13

avg
12.50

avg	count
12.50	2

Aggregation

<u>sid</u>	sname	rating	age
28	yuppy	7	15
31	debby	7	10
22	conny	5	10
58	lilly	10	13

- “Give the names of the skaters with the highest rankings”

```
SELECT sname
FROM skaters
WHERE rating = (SELECT MAX(rating)
                FROM skaters)
```

exactly 1 output
tuple

gives last that works

sname
lilly

(Note also the = in the **where** clause. We can use direct comparison (=, <, ...) when it is assured that the relation resulting from the subquery has only one tuple.)

- “Give the name of the skater that is the first in the alphabet”

```
SELECT min(sname)
FROM Skaters S
```

max is last alphabetically

min
conny

Wrong Aggregations

- “Give the names of the skaters with the highest rankings”
not equal # of tuples
SELECT sname, MAX(rating)
FROM skaters
- Does not work!
- Max is one value for ALL tuples, sname is one value for each tuple

Grouping

apply aggregation
per group

- So far, we have applied aggregate operators to all (qualifying) tuples. Sometimes, we want to apply them to each of several groups of tuples.
- Example: “Find the average age of the skaters in each rating level”
 - In general, we don’t know how many rating levels exist, and what the rating values for these levels are.
 - Suppose, we know that the rating levels go from 1 to 10; then we can write 10 queries that look like this:

```
For i = 1, 2, ... 10
SELECT AVG(age)
FROM skaters
WHERE rating = i
```

<u>sid</u>	sname	rating	age
28	yuppy	7	15
31	debby	7	10
22	conny	10	11
58	lilly	10	13

avg

12.0

Grouping

- Grouping does this with one query

```
SELECT      AVG (age) , MIN (age)
FROM skaters
GROUP BY rating
```

<u>sid</u>	sname	rating	age
28	yuppy	7	15
31	debby	7	10
22	conny	10	10
58	lilly	10	13

avg	min
12.5	10
11.5	10

for each group
get avg and
min

Queries with GROUP BY

```
SELECT target-list  
FROM relation list  
WHERE qualification  
GROUP BY grouping list
```

- A group is defined as a set of tuples that have the same value for all attributes in the grouping list
- One answer tuple is generated per group.
- The target-list contains aggregation terms and/or *attributes*
- Allowed attributes:
 - Subset of the grouping list
 - Since each answer tuple corresponds to one group, we can only depict attributes, for which all tuples in the group have the same value

- Example:

```
SELECT rating, MIN(age)  
FROM skaters  
GROUP BY rating
```

} 1 rating value per group
so, ok

Queries with GROUP BY

```
SELECT target-list  
FROM relation list  
WHERE qualification  
GROUP BY grouping list  
(ORDER BY...)
```

Example:

```
SELECT rating, MIN(age)  
FROM Skaters  
GROUP BY rating  
ORDER BY rating
```

every group
has 1 rating
to be returned,
so ok

<u>sid</u>	sname	rating	age
28	yuppy	7	15
31	debby	7	10
22	conny	10	10
58	lilly	10	13

rating	min
7	10
10	10

grouped then find min

Evaluation

```
SELECT target-list  
FROM   relation list  
(WHERE   qualification)  
GROUP BY grouping list  
(ORDER BY attributes from target-list)
```

- Conversion to Relational Algebra
 - Compute the cross-product of relations in **FROM** clause, consider only tuples that fulfill the qualification in **WHERE** clause, project on fields that are needed (in **SELECT** or **GROUP BY**)
 - Partition the remaining tuples into groups by the value of attributes in grouping-list
 - Return all attributes in the **SELECT** clause (must also be in the group list) plus the calculated aggregation terms per group.
 - Return in order if requested

SELECT lists with aggregation

- If any aggregation is used, then each element in the attribute list of the **SELECT** clause must either be aggregated or appear in a group-by clause

```
SELECT rating, MIN(age)
FROM Skater
GROUP BY rating
```

- Wrong way to find the name of the oldest skaters

```
SELECT sname, MAX(age)
FROM Skaters
```

- Correct way to find the names of the oldest skaters

```
SELECT sname, age
FROM skaters
WHERE age = (SELECT MAX(age)
              FROM skaters)
```

} otherwise use
subquery

HAVING CLAUSE

- **HAVING** clauses are selections on groups, just as **WHERE** clauses are selections on tuples
- Example: “For each rating, find the minimum age of the skaters with this rating. Only consider rating levels with at least two skaters

```
SELECT rating, MIN(age)
FROM Skaters
GROUP BY rating
HAVING COUNT(*) >= 2
```

- Example 2: “For each rating > 5, find the average age of the skaters with this rating. Only consider rating levels where there are at least two skaters

```
SELECT rating, AVG(age)
FROM Skaters
WHERE rating > 5 ← first
GROUP BY rating ← then
HAVING COUNT(*) >= 2
```

only return groups
that meet the
criteria in having
clause

<u>sid</u>	sname	rating	age
1	A	9	18
2	B	1	20
3	C	6	12
4	D	9	18
5	E	1	10
6	F	8	16
7	G	8	8

```

SELECT rating, avg(age)
  FROM Skaters
 WHERE rating > 5
 GROUP BY rating
HAVING COUNT(*) >= 2

```

- Select upon **WHERE** and project to necessary attributes
- Partition by **GROUP** and check whether they fulfill **HAVING**

rating	age	rating	age
9	18	9	18
6	12	8	13
9	18		
8	16		
8	10		

= 7 groups

<u>sid</u>	sname	rating	age
1	A	9	18
2	B	1	20
3	C	6	12
4	D	9	18
5	E	1	10
6	F	8	16
7	G	8	8

```

SELECT rating, age, count(*)
  FROM Skaters
 WHERE rating > 5
 GROUP BY rating, age
HAVING COUNT(*) >= 2

```

need
same
value
for
both

rating	age	count
9	18	2

'where' is executed first

Evaluation

```
SELECT rating, avg(age)
FROM Skaters
WHERE rating > 5
GROUP BY rating
HAVING COUNT(*) >= 2
```

```
SELECT target-list
FROM relation list
WHERE qualification
GROUP BY grouping list
HAVING group-qualification
```

- Conversion to Relational Algebra

- Compute the cross-product of relations in **FROM** clause, consider only tuples that fulfill the qualification in **WHERE** clause, project on fields that are needed (in **SELECT** or **GROUP BY**)
- NOTE: the **WHERE** clause can contain any attributes of the relations in the relation list

```
SELECT rating, MIN(age)
FROM Skaters
WHERE sname LIKE 'A%'
GROUP BY rating
HAVING COUNT(*) >= 2
```

- Partition the remaining tuples into groups by the value of attributes in grouping-list
- For each group, the group qualification is then applied selecting only those groups that fulfill the qualification. Expressions in group-qualification must have a single value per group. Hence, for each attribute in the group qualification, either
 - the attribute also appears in the grouping list
 - or it is argument of an aggregation

Example II

- Grouping over a join of two relations.
- *For each local competition, find the number of participants*


```
SELECT  C.cid, COUNT (*) AS scount
FROM    Competition C, Participates P
WHERE   C.cid=P.cid AND c.type='local'
GROUP BY C.cid
```


Example III

- Find those ratings for which the average age is the minimum over all ratings

- Aggregate operations cannot be nested! **WRONG**

```
SELECT S.rating
FROM   Skaters S
WHERE  S.age = (SELECT MIN (AVG (S2.age))
                FROM Skaters S2)
```



- Use views

complex having clause

```
select S.rating, Avg(S.age)
from Skaters S
group by S.rating
having Avg(S.age) <= All
      (select Avg(S2.age)
       from Skaters S2)
```

Similar to renaming in relational algebra

group by skating

Views

- A view is just an unmaterialized relation: we store a definition rather than a set of tuples.

```
CREATE VIEW ActiveSkaters (sid, sname)
AS SELECT DISTINCT S.sid, S.sname
FROM Skaters S, Participates P
WHERE S.sid = P.sid
```

- Views can be used to present necessary information (or a summary), while hiding details in underlying relation(s).

☆ Given ActiveSkaters, we know the names of the skaters who have participated competition but not the age of the skaters (may be uninteresting for the users of ActiveSkaters).

make base tables smaller

Views

- Use a **view as intermediate relation** (rename in rel.algebra)
- *Find those ratings for which the average age is the minimum over all ratings*

```
CREATE VIEW Temp (rating, avgage)
AS SELECT rating, AVG (age) AS avgage
FROM skaters
GROUP BY rating)
```

```
SELECT rating, avgage
FROM Temp
WHERE avgage = (SELECT MIN (avgage)
                FROM Temp)
```

} view
definition
query

} query that
uses view
as a table

Views (contd)

- Views can be treated as if they were materialized relations
- The system translates a **SELECT** on a view into **SELECTS** on the materialized relations
- Modifications are problematic
- Views can be dropped using the **DROP VIEW** command
 - How to handle **DROP TABLE** if there's a view on the table?
 - **DROP TABLE** command has options to let the user specify this.

NULL Values

- Meaning of a NULL value
 - Unknown/missing
 - Inapplicable (e.g., no spouse's name)
- Comparing NULLs to values
 - E.g., how to evaluate condition $\text{rating} > 7$ if tuple has a NULL in rating?
 - When we compare a NULL value and any other value (including NULL) using a comparison operator like $>$ or $=$, the result is “unknown”.
 - If we want to check whether a value is NULL, SQL provides the special comparison operator **IS NULL**
- Arithmetic Operations (*, +, etc):
 - When at least one operand has a NULL value (the other operands can have any value including NULL) then the result is NULL (consequence $0 * \text{NULL} = \text{NULL} !$)
 - We cannot use NULL as an operand (e.g., $\text{rating} < \text{NULL}$).

NULL Values (contd.)

- 3-valued logic necessary: **true**, **false**, **unknown**

- NOT A

A	NOT A
true	false
false	true
unknown	unknown

- A AND B

A \ B	true	false	unknown
true	true	false	unknown
false	false	false	false
unknown	unknown	false	unknown

- A OR B

A \ B	true	false	unknown
true	true	true	true
false	true	false	unknown
unknown	true	unknown	unknown

Query evaluation considering NULL values

- Evaluation in SQL
 - The qualification in the **WHERE** clause eliminates rows for which the qualification does not evaluate **true** (i.e., rows that evaluate to **false** or **unknown** are eliminated)
 - SQL defines that rows are duplicates if corresponding columns are either equal or both contain **NULL** (in contrast to the usual on previous slide where the comparison of the **NULLs** results in **unknown**)
 - **COUNT (*)** handles **NULLs** like other values, i.e., they are counted
 - All other aggregate operations simply discard **NULL** values

Inner Join (Default)

Dangling tuples: No match in the other relation → no output

```
SELECT *
```

```
FROM skaters s INNER JOIN participates p
```

```
ON s.sid = p.sid
```

```
skaters
```

Optional Keyword

normal join

<u>sid</u>	sname	rating	age
28	yuppy	9	15
31	debby	7	10
22	conny	5	10
58	lilly	10	13

participates

<u>sid</u>	<u>cid</u>	rank
31	101	2
58	103	7
58	101	7

S.Sid	sname	rating	age	p.sid	cid	rank
31	debby	7	10	31	101	2
58	lilly	10	13	58	103	7
58	lilly	10	13	58	101	7

Outer Join

Dangling tuples: No match in the other relation → One dummy record

SELECT *

FROM skaters s LEFT OUTER JOIN participates p

WHERE s.sid = p.sid

every tuple in
left table
appears in
output

sid	sname	rating	age
28	yuppy	9	15
31	debby	7	10
22	conny	5	10
58	lilly	10	13

sid	cid	rank
31	101	2
58	103	7
58	101	7

S.Sid	sname	rating	age	p.sid	cid	rank
28	yuppy	9	15	NULL	NULL	NULL
31	debby	7	10	31	101	2
22	conny	5	10	NULL	NULL	NULL
58	lilly	10	13	58	103	7
58	lilly	10	13	58	101	7

no matching
tuple
but
included

Outer Join Types

SELECT *
FROM A **LEFT OUTER JOIN** B
WHERE A.att1 = B.att2

Pad dangling
tuples from A

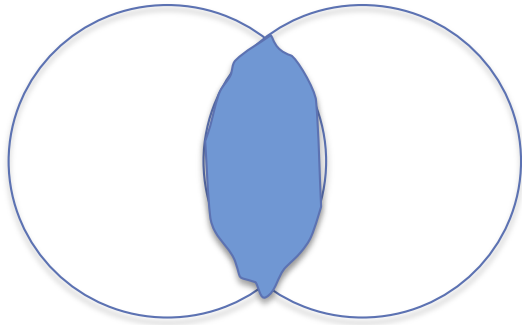
SELECT *
FROM A **RIGHT OUTER JOIN** B
WHERE A.att1 = B.att2

Pad dangling
tuples from B

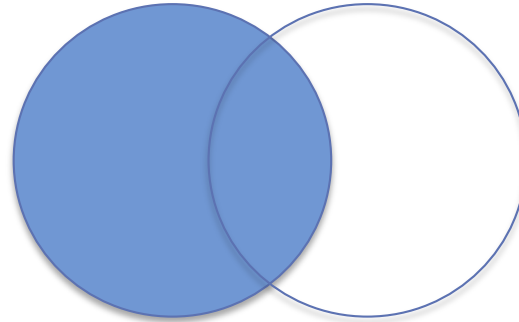
SELECT *
FROM A **FULL OUTER JOIN** B
WHERE A.att1 = B.att2

Pad dangling
tuples from A and
B

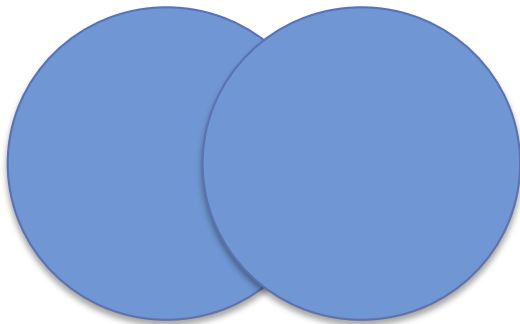
Visualization



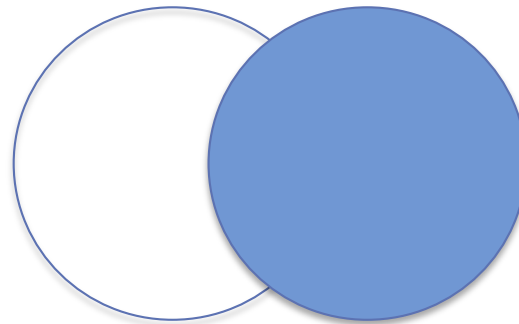
A INNER JOIN B



A LEFT OUTER JOIN



A FULL OUTER JOIN



A RIGHT OUTER JOIN

Foreign Key Constraints

- If B has NOT NULL foreign key referencing A
 - Every tuple in B joins with exactly one tuple in A
 - $A \text{ INNER JOIN } B = A \text{ RIGHT OUTER JOIN } B$
 - As no dangling tuples

WITH CLAUSE

Functions like the concept of a view for the scope of this SQL statement

WITH *new table name* partinfo(sid, sname, age, rank, cid) **AS**

(

SELECT S.sid, S.sname, S.age, P.rank, P.cid

FROM skaters S

INNER JOIN participates P

ON S.sid = P.sid

)

SELECT sid, sname, cid

FROM partinfo

WHERE age > 7;

attributes

Can be any complex SQL Select query, with joins, aggregations, etc.

Can be any complex SQL Select query, with joins, aggregations, etc., including with other tables.

WITH A(...) AS
(SELECT ...)
, B(...) AS
(SELECT ...)
SELECT ...;

Possible to have multiple aliases defined

Makes temp table

temp table definition in from clause

Derived Tables

```
SELECT sid, sname, cid
```

```
FROM
```

```
(
```

```
    SELECT S.sid, S.sname, S.age, P.rank, P.cid
```

```
    FROM skaters S
```

```
        INNER JOIN participates P
```

```
        ON S.sid = P.sid
```

```
)partinfo name
```

```
WHERE age > 7;
```

Can be any complex SQL Select query, with joins, aggregations, etc.

Functions like the concept of a view for the scope of this SQL statement

The outer query can be any complex SQL Select query, with joins, aggregations, etc., including with other tables.

```
SELECT ...  
FROM (SELECT ...)A  
     , (SELECT ...)B  
WHERE ...
```

Possible to have multiple aliases defined

With temp (rating, avgage) AS

(select rating, Avg(age)

from skaters

Group by rating)

select rating, avgage

from temp

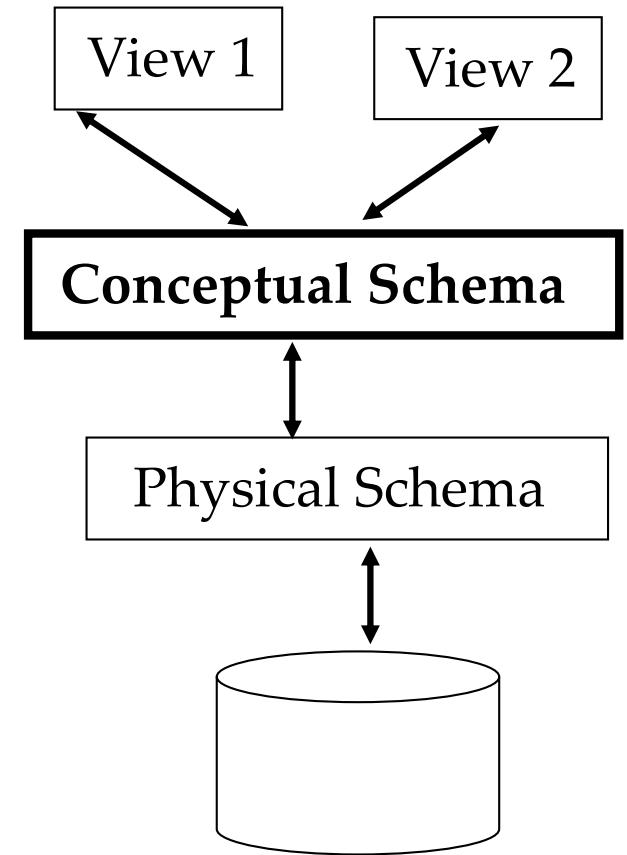
where avgage = (select min (avgage)

from temp);

ratings for which avgage is min

Levels of Abstraction

- ❑ Single conceptual (logical) schema defines logical structure
 - ☆ Conceptual database design
- ❑ Physical schema describes the files and indexes used
 - ☆ Physical database design
- ❑ Different views describe how users see the data (also referred to as external schema)
 - ☆ generated on demand from the real data
- ❑ Physical data independence: the conceptual schema protects from changes in the physical structure of data
- ❑ Logical data independence: external schema protects from changes in conceptual schema of data



DB Modifications: insert/delete/update

- Insert values for all attributes in the order attributes were declared or values for only some attributes
 - `INSERT INTO Skaters VALUES (68, 'Jacky', 10, 10)` *all*
 - `INSERT INTO Skaters (sid, name) VALUES (68, 'Jacky')` *some*
- Insert the result of a query
 - `ActiveSkaters (sid, name)`
 - `INSERT INTO ActiveSkaters (
SELECT Skaters.sid Skaters.name
FROM Skaters, Participates
WHERE Skaters.sid = Participates.sid)`

DB Modifications: insert/delete/update

- Delete some or all tuples of a relation
 - `DELETE FROM Competitions WHERE cid = 103` *some*
 - `DELETE FROM Competitions` *all*
- Update some of the attributes of some of the tuples
 - `UPDATE Skaters`
 - `SET ranking = 10, age = age + 1`
 - `WHERE name = 'debby' OR name = 'lilly'`
- SQL2 semantics: all conditions in a modification statement must be evaluated by the system *BEFORE* any modifications occur.