

SQL: Part I

CS348 Spring 2024

Instructor: Sujaya Maiyya

Sections: **002 and 003 only**

SQL

- SQL: **Structured Query Language**
 - Pronounced “S-Q-L” or “sequel”
 - The standard query language supported by most DBMS
 - Introduced in 1970s and standardized by ANSI since 1986

SQL

- **Data-definition language (DDL):** define/modify schemas, delete relations
- **Data-manipulation language (DML):** query information, and insert/delete/modify tuples
- **Integrity constraints:** specify constraints that the data stored in the database must satisfy
- Intermediate/Advanced topics: **(next week)**
 - E.g., triggers, views, indexes, programming, recursive queries

this
week

DDL

User (uid int, name string, age int, pop float)
Group (gid string, name string)
Member (uid int, gid string)

- **CREATE TABLE** *table_name*
(..., *column_name column_type*, ...);

```
CREATE TABLE User(uid INT, name VARCHAR(30), age INT, pop DECIMAL(3,2));  
CREATE TABLE Group (gid CHAR(10), name VARCHAR(100));  
CREATE TABLE Member (uid INT, gid CHAR(10));
```

- **DROP TABLE** *table_name*;

```
DROP TABLE User;  
DROP TABLE Group;  
DROP TABLE Member;
```

Drastic action:
deletes ALL info
about the table, not
just the contents

-- everything from -- to the end of line is ignored.
-- SQL is insensitive to white space.
-- SQL is insensitive to case (e.g., ...CREATE... is equivalent to ...create...).

Basic queries for DML: SFW statement

- **SELECT** A_1, A_2, \dots, A_n
FROM R_1, R_2, \dots, R_m
WHERE *condition*;
- Also called an SPJ (select-project-join) query
- Corresponds to (**but not really equivalent to**) relational algebra query:

$$\pi_{A_1, A_2, \dots, A_n} \left(\sigma_{\text{condition}} (R_1 \times R_2 \times \dots \times R_m) \right)$$

Examples

User (uid int, name string, age int, pop float)
Group (gid string, name string)
Member (uid int, gid string)

- List all rows in the User table

```
SELECT * FROM User;
```

- * is a short hand for “all columns”

- List name of users under 18 (selection, projection)

```
SELECT name FROM User where age <18;
```

- When was Lisa born?

```
SELECT 2024-age FROM User where name = 'Lisa';
```

- SELECT list can contain expressions
- String literals (case sensitive) are enclosed in quotes

Example: join

User (uid int, name string, age int, pop float)
Group (gid string, name string)
Member (uid int, gid string)

- List IDs and names of groups with a user whose name contains “Simpson”

```
SELECT Group.gid, Group.name
      FROM User, Member, Group
     WHERE User.uid = Member.uid
           AND Member.gid = Group.gid
           AND ...;
```

Example: join

User (uid int, name string, age int, pop float)
Group (gid string, name string)
Member (uid int, gid string)

- List ID's and names of groups with a user whose name **contains** "Simpson"

```
SELECT Group.gid, Group.name
      FROM User, Member, Group
      WHERE User.uid = Member.uid
            AND Member.gid = Group.gid
            AND User.name LIKE '%Simpson%';
```

- **LIKE** matches a string against a pattern
 - **%** matches any sequence characters
- Okay to omit *table_name* in *table_name.column_name* if *column_name* is unique

Example: rename

User (uid int, name string, age int, pop float)
Group (gid string, name string)
Member (uid int, gid string)

- IDs of all pairs of users that belong to one group
 - Relational algebra query:

$$\pi_{m_1.uid, m_2.uid} (\rho_{m_1} Member \bowtie_{m_1.gid=m_2.gid \wedge m_1.uid > m_2.uid} \rho_{m_2} Member)$$

- SQL (not exactly due to duplicates):

```
SELECT m1.uid AS uid1, m2.uid AS uid2
       FROM Member AS m1, Member AS m2
       WHERE m1.gid = m2.gid
             AND m1.uid > m2.uid;
```

- **AS** keyword is completely optional

A more complicated example

- Names of all groups that Lisa and Ralph are both in

Tip: Write the FROM clause first, then WHERE, and then SELECT

User (uid int, name string, age int, pop float)
Group (gid string, name string)
Member (uid int, gid string)

A more complicated example

- Names of all **groups that Lisa** and Ralph are both in

```
SELECT g.name
  FROM User u1, ..., Member m1, ...
 WHERE u1.name = 'Lisa' AND ...
        AND u1.uid = m1.uid AND ...
        AND ...;
```

User (uid int, name string, age int, pop float)
Group (gid string, name string)
Member (uid int, gid string)

A more complicated example

- Names of all **groups that** Lisa and **Ralph** are both in

```
SELECT g.name
  FROM User u1, User u2, Member m1, Member m2, ...
 WHERE u1.name = 'Lisa' AND u2.name = 'Ralph'
       AND u1.uid = m1.uid AND u2.uid=m2.uid
       AND ...;
```

User (uid int, name string, age int, pop float)
Group (gid string, name string)
Member (uid int, gid string)

A more complicated example

- Names of all groups that Lisa and Ralph are both in

```
SELECT g.name
FROM User u1, User u2, Member m1, Member m2, Group g
WHERE u1.name = 'Lisa' AND u2.name = 'Ralph'
      AND u1.uid = m1.uid AND u2.uid=m2.uid
      AND m1.gid = g.gid AND m2.gid = g.gid;
```

User (uid int, name string, age int, pop float)
Group (gid string, name string)
Member (uid int, gid string)

Why SFW statements?

- Many queries can be written using only **selection, projection, and cross product (or join)**
- These queries can be written in a canonical form which is captured by SFW:

$$\pi_L \left(\sigma_p (R_1 \times \cdots \times R_m) \right)$$

- E.g.: $\pi_{R.A,S.B} (R \bowtie_{p_1} S) \bowtie_{p_2} (\pi_{T.C} \sigma_{p_3} T)$ can be written as
 $= \pi_{R.A,S.B,T.C} \sigma_{p_1 \wedge p_2 \wedge p_3} (R \times S \times T)$

Set versus bag

User

| uid | name | age | pop |
|-----|----------|-----|-----|
| 142 | Bart | 10 | 0.9 |
| 123 | Milhouse | 10 | 0.2 |
| 857 | Lisa | 8 | 0.7 |
| 456 | Ralph | 8 | 0.3 |
| ... | ... | ... | ... |

age

10

8

...

$\pi_{age} User$

Set

- No duplicates
- Relational model and algebra use set semantics

```
SELECT age  
FROM User;
```

age

10

10

8

8

...

Bag

- Duplicates allowed
- Rows in output = rows in input (w/o where clause)
- SQL uses bag semantics by default

A case for bag semantics

- Efficiency
 - Saves time of eliminating duplicates

- Which one is more useful?

$\pi_{age} User$

```
SELECT age  
FROM User;
```

- The first query just returns all possible user ages in the table
 - The second query returns the user age distribution
- Besides, SQL provides the option of set semantics with **DISTINCT** keyword

Forcing set semantics

- IDs of all pairs of users that belong to one group

```
SELECT m1.uid AS uid1, m2.uid AS uid2
       FROM Member AS m1, Member AS m2
       WHERE m1.gid = m2.gid
              AND m1.uid > m2.uid;
```

→ Say Lisa and Ralph are in both the book club and the student government, their id pairs will appear twice

- Remove duplicate (uid1, uid2) pairs from the output

```
SELECT DISTINCT m1.uid AS uid1, m2.uid AS uid2
       FROM Member AS m1, Member AS m2
       WHERE m1.gid = m2.gid;
              AND m1.uid > m2.uid;
```

Semantics of SFW

- **SELECT [DISTINCT] E_1, E_2, \dots, E_n**
FROM R_1, R_2, \dots, R_m
WHERE *condition*;
- For each t_1 in R_1 :
 For each t_2 in R_2 :
 For each t_m in R_m :
 If *condition* is true over t_1, t_2, \dots, t_m :
 Compute and output E_1, E_2, \dots, E_n as a row
 If DISTINCT is present
 Eliminate duplicate rows in output
- t_1, t_2, \dots, t_m are often called **tuple variables**

SQL set and bag operations

- Set: UNION, EXCEPT, INTERSECT
 - Exactly like set \cup , $-$, and \cap in relational algebra
 - Duplicates in input tables, if any, are first eliminated
 - Duplicates in result are also eliminated (for UNION)

| Bag1 | Bag2 |
|--------------|--------------|
| <i>fruit</i> | <i>fruit</i> |
| apple | orange |
| apple | orange |
| orange | orange |

(SELECT * FROM Bag1)
UNION
(SELECT * FROM Bag2);

| |
|--------------|
| <i>fruit</i> |
| apple |
| orange |

(SELECT * FROM Bag1)
EXCEPT
(SELECT * FROM Bag2);

| |
|--------------|
| <i>fruit</i> |
| apple |

(SELECT * FROM Bag1)
INTERSECT
(SELECT * FROM Bag2);

| |
|--------------|
| <i>fruit</i> |
| orange |

SQL set and bag operations

- Set: UNION, EXCEPT, INTERSECT
 - Exactly like set \cup , $-$, and \cap in relational algebra
- Bag: UNION ALL, EXCEPT ALL, INTERSECT ALL
 - Think of each row as having an implicit **count** (the number of times it appears in the table)

| Bag1 | Bag2 | | | | | | | | |
|--|-----------------------|-------|-------|--------|---|-------|-------|--------|--------|
| <table><thead><tr><th>fruit</th></tr></thead><tbody><tr><td>apple</td></tr><tr><td>apple</td></tr><tr><td>orange</td></tr></tbody></table> | fruit | apple | apple | orange | <table><thead><tr><th>fruit</th></tr></thead><tbody><tr><td>apple</td></tr><tr><td>orange</td></tr><tr><td>orange</td></tr></tbody></table> | fruit | apple | orange | orange |
| fruit | | | | | | | | | |
| apple | | | | | | | | | |
| apple | | | | | | | | | |
| orange | | | | | | | | | |
| fruit | | | | | | | | | |
| apple | | | | | | | | | |
| orange | | | | | | | | | |
| orange | | | | | | | | | |
| apple: 2 orange: 1 | apple: 1 orange: 2 | | | | | | | | |

```
(SELECT * FROM Bag1)  
UNION ALL  
(SELECT * FROM Bag2);
```

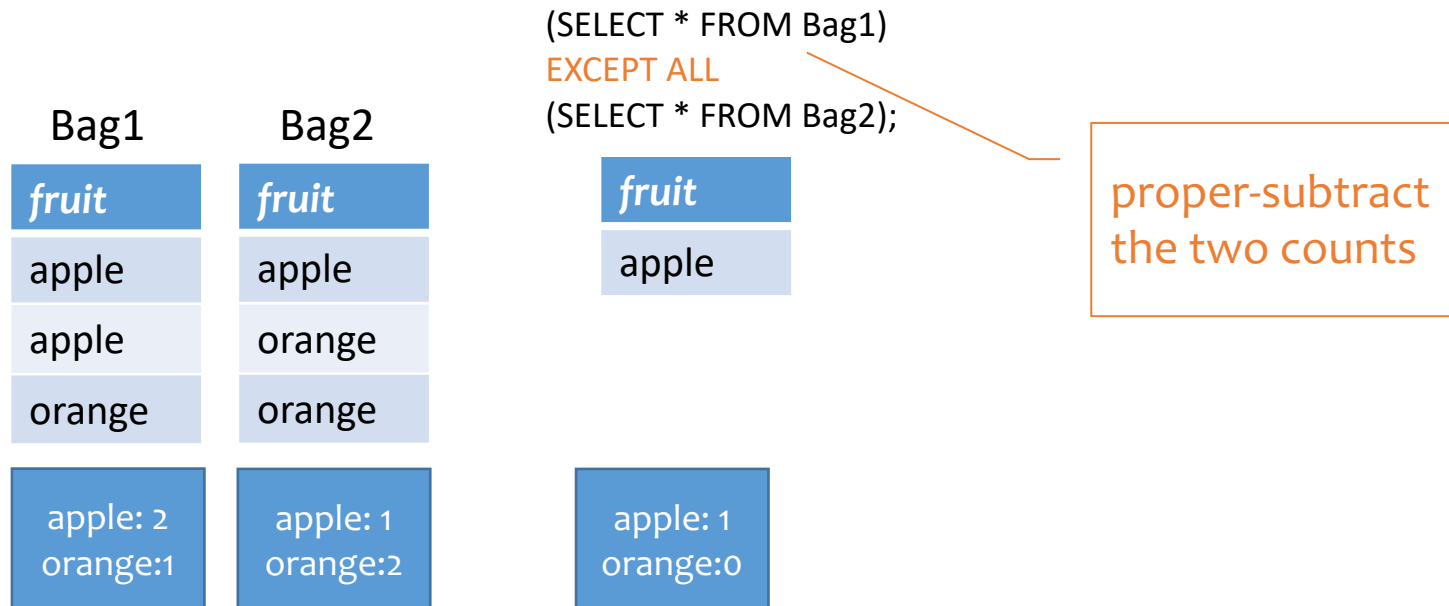
| fruit |
|--------|
| apple |
| apple |
| orange |
| apple |
| orange |
| orange |

| |
|-----------------------|
| apple: 3 orange: 3 |
|-----------------------|

sum up the counts
from two tables

SQL set and bag operations

- Set: UNION, EXCEPT, INTERSECT
 - Exactly like set \cup , $-$, and \cap in relational algebra
- Bag: UNION ALL, EXCEPT ALL, INTERSECT ALL
 - Think of each row as having an implicit **count** (the number of times it appears in the table)



SQL set and bag operations

- Set: UNION, EXCEPT, INTERSECT
 - Exactly like set \cup , $-$, and \cap in relational algebra
- Bag: UNION ALL, EXCEPT ALL, INTERSECT ALL
 - Think of each row as having an implicit **count** (the number of times it appears in the table)

| Bag1 | Bag2 | | | | | | | | |
|--|-----------------------|-------|-------|--------|---|-------|-------|--------|--------|
| <table><thead><tr><th>fruit</th></tr></thead><tbody><tr><td>apple</td></tr><tr><td>apple</td></tr><tr><td>orange</td></tr></tbody></table> | fruit | apple | apple | orange | <table><thead><tr><th>fruit</th></tr></thead><tbody><tr><td>apple</td></tr><tr><td>orange</td></tr><tr><td>orange</td></tr></tbody></table> | fruit | apple | orange | orange |
| fruit | | | | | | | | | |
| apple | | | | | | | | | |
| apple | | | | | | | | | |
| orange | | | | | | | | | |
| fruit | | | | | | | | | |
| apple | | | | | | | | | |
| orange | | | | | | | | | |
| orange | | | | | | | | | |
| apple: 2 orange: 1 | apple: 1 orange: 2 | | | | | | | | |

```
(SELECT * FROM Bag1)  
INTERSECT ALL  
(SELECT * FROM Bag2);
```

| fruit |
|--------|
| apple |
| orange |

| |
|-----------------------|
| apple: 1 orange: 1 |
|-----------------------|

take the
minimum of the
two counts

Set versus bag operations

Poke (uid1, uid2, timestamp)

- uid1 poked uid2 at timestamp

Question: How do these two queries differ?

Q1:

```
(SELECT uid1 FROM Poke)
```

EXCEPT

```
(SELECT uid2 FROM Poke);
```

Q2:

```
(SELECT uid1 FROM Poke)
```

EXCEPT ALL

```
(SELECT uid2 FROM Poke);
```

Set versus bag operations

Poke (uid1, uid2, timestamp)

- uid1 poked uid2 at timestamp

Question: How do these two queries differ?

Q1:

```
(SELECT uid1 FROM Poke)
```

EXCEPT

```
(SELECT uid2 FROM Poke);
```

Users who poked others but
never got poked by others

Q2:

```
(SELECT uid1 FROM Poke)
```

EXCEPT ALL

```
(SELECT uid2 FROM Poke);
```

Users who poked others
more than others poked them

SQL features covered so far

- Query
 - SELECT-FROM-WHERE statements
 - Set/bag (DISTINCT, UNION/EXCEPT/INTERSECT (ALL))

👉 Next: how to **nest SQL queries**

- Subqueries (table, scalar, IN, EXISTS, ALL, ANY)
- Aggregation and grouping (GROUP BY, HAVING)
- Ordering (ORDER)
- Joins