Intro to Database Systems (15-445/645) **Modern SQL** Carnegie Mellon **SPRING** Charlie 2023 Garrod **Jniversity**

ADMINISTRIVIA

Project 0 due Sunday. Quick adaptation to C++ is a prerequisite for this course!

Homework 1 released today, due Friday Feb 3rd.

LAST CLASS

We introduced the Relational Model as the superior data model for databases.

We then showed how Relational Algebra is the building blocks that will allow us to query and modify a relational database.

ASIDE: OTHER DATA MODELS

Relational
Key/Value
GraphImmongoDB
mongoDBDocument / ObjectLeading AlternativeWide-Column / Column-family
Array / Matrix / VectorsImmongoDB
elasticHierarchical
NetworkImmongoDB
elasticMulti-ValueImmongoDB
mongoDB



RELATIONAL MODEL: QUERIES

The relational model is independent of any query language implementation.

SQL is the *de facto* standard (many dialects).

for line in file.readlines():
 record = parse(line)
 if record[0] == "GZA":
 print(int(record[1]))

SELECT year FROM artists
WHERE name = 'GZA';

SQL

In 1971, IBM created language called <u>SQUA</u>

IBM then created "SH System R prototype I \rightarrow <u>S</u>tructured <u>English</u> <u>Q</u>u

Q2. Find the average salary of employees in the Shoe Department. AVG (EMP' ('SHOE'))

Mappings may be composed by applying one mapping to the result of another, as illustrated by Q3.

Q3. Find those items sold by departments on the second floor. SALES DEPT DEPT LOC FLOOR (2) ITEM

The floor '2' is first mapped to the departments located there, and then to the items which they sell. The range of the inner mapping must be compatible with the domain of the outer mapping, but they need not be identical, as illustrated by Q4. IBM releases commerciar or

→ System/38 (1979), SQL/DS (1981), and DB2 (1983).

Security CMU.DB 15-445/645 (Spring 2023



The minimum language syntax a system needs to say that it supports SQL is <u>SQL-92</u>.

RELATIONAL LANGUAGES

Data Manipulation Language (DML) Data Definition Language (DDL) Data Control Language (DCL)

Also includes:

- \rightarrow View definition
- → Integrity & Referential Constraints
- \rightarrow Transactions

Important: SQL is based on **multisets** (a.k.a. **bags**, with duplicates), not **sets** (no duplicates).

TODAY'S AGENDA

10

Aggregations + Group By String / Date / Time Operations Output Control + Redirection Nested Queries Window Functions Common Table Expressions

EXAMPLE DATABASE

student(sid,name,login,gpa)

| sid | name | login | age | gpa |
|-------|--------|------------|-----|-----|
| 53666 | Kanye | kanye@cs | 44 | 4.0 |
| 53688 | Bieber | jbieber@cs | 27 | 3.9 |
| 53655 | Тирас | shakur@cs | 25 | 3.5 |

course(cid,name)

| С | id | name |
|---|-------|-----------------------------|
| 1 | 5-445 | Database Systems |
| 1 | 5-721 | Advanced Database Systems |
| 1 | 5-826 | Data Mining |
| 1 | 5-799 | Special Topics in Databases |

enrolled(sid, cid, grade)

11

| sid | cid | grade |
|-------|--------|-------|
| 53666 | 15-445 | С |
| 53688 | 15-721 | А |
| 53688 | 15-826 | В |
| 53655 | 15-445 | В |
| 53666 | 15-721 | С |

AGGREGATES

Functions that return a single value from a bag of tuples:

- \rightarrow AVG(col) \rightarrow Return the average col value.
- \rightarrow MIN(col) \rightarrow Return minimum col value.
- \rightarrow MAX(col) \rightarrow Return maximum col value.
- \rightarrow SUM(col) \rightarrow Return sum of values in col.
- \rightarrow COUNT(col) \rightarrow Return # of values for col.

AGGREGATES

13

Aggregate functions can (almost) only be used in the **SELECT** output list.

Get # of students with a "@cs" login:



MULTIPLE AGGREGATES

Get the number of students and their average GPA that have a "@,cs" login.

| | | | AVG(gpa) | COUNT (| sid) |
|---|--------|--------------------------|----------|---------|------|
| _ | SELECT | AVG(gpa), COUNT(sid) | 3.8 | 3 | |
| V | FROM | student WHERE login LIKE | '%@cs' | | |

DISTINCT AGGREGATES

COUNT, SUM, AVG support DISTINCT

Get the number of unique students that have an "@cs" login.

COUNT(DISTINCT login)SELECT COUNT(DISTINCT login)3FROM student WHERE login LIKE'%@cs'

AGGREGATES

Output of other columns outside of an aggregate is undefined.

Get the average GPA of students enrolled in each course.

ECMU-DB 15-445/645 (Spring 2023

GROUP BY

Project tuples into subsets and calculate aggregates against each subset.

| SELECT | AVG(s.gpa), e.cid |
|--------|---------------------------------|
| FROM | enrolled AS e JOIN student AS s |
| ON | e.sid = s.sid |
| GROUP | BY e.cid |

17

| e.sid | s.sid | s.gpa | e.cid |
|-------|-------|-------|--------|
| 53435 | 53435 | 2.25 | 15-721 |
| 53439 | 53439 | 2.70 | 15-721 |
| 56023 | 56023 | 2.75 | 15-826 |
| 59439 | 59439 | 3.90 | 15-826 |
| 53961 | 53961 | 3.50 | 15-826 |
| 58345 | 58345 | 1.89 | 15-445 |

| AVG(s.gpa) | e.cid |
|------------|--------|
| 2.46 | 15-721 |
| 3.39 | 15-826 |
| 1.89 | 15-445 |

GROUP BY

Non-aggregated values in SELECT output clause must appear in GROUP BY clause.

SELECT AVG(s.gpa), e.cid, s. ne
FROM enrolled AS e JOIN stude t AS s
ON e.sid = s.sid
GROUP BY e.cid, s.name

HAVING

Filters results based on aggregation computation. Like a WHERE clause for a GROUP BY

SELECT AVG(s.gpa) AS avg_gpa, e.cid
FROM enrolled AS e, student AS s
WHERE e.sid = s.sid
GROUP BY e.cid
HAVING AVG(s.gpa) > 3.9;

| AVG(s.gpa) | e.cid | | |
|------------|--------|----------|--------|
| 3.75 | 15-415 | avg_gpa | e.cid |
| 3.950000 | 15-721 | 3.950000 | 15-721 |
| 3.900000 | 15-826 | | |

| | String Case | String Quotes | |
|------------|-----------------|--------------------------------|--|
| SQL-92 | Sensitive | Single Only | |
| Postgres | Sensitive | Single Only | |
| MySQL | Insensitive | Single/Double | |
| SQLite | Sensitive | Single/Double | |
| MSSQL | Sensitive | Single Only | |
| Oracle | Sensitive | Single Only | |
| WHERE UPPI | ER(name) = UPPE | R('KaNyE') <mark>SQL-92</mark> | |
| WHERE name | e = "KaNyE" | MySQL | |
| | | | |

LIKE is used for string matching.
String-matching operators
→ '%' Matches any substring (including empty strings).
→ '_' Match any one character

SELECT * FROM enrolled AS e
WHERE e.cid LIKE '15-%'

SELECT * FROM student AS s
WHERE s.login LIKE '%@c_'

ECMU-DB 15-445/645 (Spring 2023

SQL-92 defines string functions. → Many DBMSs also have their own unique functions Can be used in either output and predicates:

SELECT SUBSTRING(name,1,5) AS abbrv_name
FROM student WHERE sid = 53688

```
SELECT * FROM student AS s
WHERE UPPER(s.name) LIKE 'KAN%'
```

SQL standard says to use || operator to concatenate two or more strings together.

SELECT name FROM studentSQL-92WHERE login = LOWER(name) || '@cs'SELECT name FROM studentMSSQL

WHERE login = LOWER(name) + '@cs'

SELECT name FROM student MySQL
WHERE login = CONCAT(LOWER(name), '@cs')

DATE/TIME OPERATIONS

Operations to manipulate and modify DATE/TIME attributes.

Can be used in both output and predicates.

Support/syntax varies wildly...

OUTPUT REDIRECTION

25

Store query results in another table:

- \rightarrow Table must not already be defined.
- \rightarrow Table will have the same # of columns with the same types as the input.

| SELECT | DISTINCT cid INTO CourseIds SQL-92 | |
|--------|------------------------------------|----------|
| FROM | SELECT DISTINCT cid | Postgres |
| | INTO TEMPORARY CourseIds | |
| CREATE | <pre>FROM enrolled;</pre> | |
| SELE | CT DISTINCT cid FROM enrolled); | |

OUTPUT REDIRECTION

Insert tuples from query into another table:

- \rightarrow Inner **SELECT** must generate the same columns as the target table.
- → DBMSs have different options/syntax on what to do with integrity violations (e.g., invalid duplicates).

INSERT INTO CourseIds SQL-92
(SELECT DISTINCT cid FROM enrolled);

ECMU-DB 15-445/645 (Spring 2023

OUTPUT CONTROL

ORDER BY <column*> [ASC|DESC]

ECMU·DB 15-445/645 (Spring 2023) → Order the output tuples by the values in one or more of their columns.





OUTPUT CONTROL

LIMIT <count> [offset]

 \rightarrow Limit the # of tuples returned in output.

 \rightarrow Can set an offset to return a "range"

| SELECT WHERE | sid, name FROM student login LIKE '%@cs' | |
|-----------------|---|-------------|
| LIMIT | 10 SELECT TOP 10 sid, name FROM stu | udent MSSQL |
| SELECT | sid, WHERE login LIKE '%@cs' | |
| WHERE | login LIKE '%@cs' | |
| LIMIT | 10 OFFSET 20 | |
| | | |

28

Queries containing other queries. They are often difficult to optimize.

Inner queries can appear (almost) anywhere in query.

Outer Query → SELECT name FROM student WHERE sid IN (SELECT sid FROM enrolled) ← Inner Query



 $ALL \rightarrow$ Must satisfy expression for all rows in the sub-query.

 $ANY \rightarrow Must satisfy expression for at least one row in the sub-query.$

 $IN \rightarrow Equivalent to '=ANY()'$.

EXISTS \rightarrow At least one row is returned without comparing it to an attribute in outer query.

Get the names of students in '15-445'

```
SELECT name FROM student
WHERE sid = ANY(
   SELECT sid FROM enrolled
   WHERE cid = '15-445'
```

Find student record with the highest id that is enrolled in at least one course.



This won't work in SQL-92. It runs in SQLite, but not Postgres or MySQL (v8 with strict mode).

34

Find student record with the highest id that is enrolled in at least one course.



15-445/645 (Spring 2023

Security CMU.DB

Find all courses that have no students enrolled in it.

 SELECT * FROM course

 WHERE NOT EXISTS(

 SELECT * FROM enrolled

 WHERE course.cid = enrolled.cid

 15-445

 Database Systems

 53666

 15-445

| 15-115 Databasa Systems | | | | | - | |
|-------------------------|----------------|----------------------|---|----------|---------|---|
| 15 445 | Dalabase Syste | 1115 | | Facaa | 1 - 701 | |
| 15-721 | | | | | | A |
| 13 721 | | name | | | | R |
| 15-826 | Data | | | | | |
| 15-620 | 15-799 | Special Topics in Da | | atahases | | Ь |
| 15-700 | | | | LUDUSCS | | D |
| 10-199 | special topics | III Databases | | 52666 | 15 701 | |
| / | | | - | 0000 | 13-721 | |

Performs a "sliding" calculation across a set of tuples that are related.

Like an aggregation but tuples are not grouped into a single output tuples.

> How to "slice" up data Can also sort

SELECT ... FUNC-NAME(...) OVER (...)

FROM tableName

Aggregation Functions Special Functions

ECMU·DB 15-445/645 (Spring 2023

Aggregation functions:
→ Anything that we discussed earlier
Special window functions:
→ ROW_NUMBER() → # of the current row
→ RANK() → Order position of the current row.

| sid | cid | grade | row_num |
|-------|--------|-------|---------|
| 53666 | 15-445 | С | 1 |
| 53688 | 15-721 | А | 2 |
| 53688 | 15-826 | В | 3 |
| 53655 | 15-445 | В | 4 |
| 53666 | 15-721 | С | 5 |

SELECT *, ROW_NUMBER() OVER () AS row_num
FROM enrolled

The OVER keyword specifies how to group together tuples when computing the window function. Use PARTITION BY to specify group.

| cid | sid | row_number | |
|--------|-------|------------|--|
| 15-445 | 53666 | 1 | |
| 15-445 | 53655 | 2 | |
| 15-721 | 53688 | 1 | |
| 15-721 | 53666 | 2 | |
| 15-826 | 53688 | 1 | |

SELECT cid, sid, ROW_NUMBER() OVER (PARTITION BY cid) FROM enrolled ORDER BY cid

You can also include an **ORDER BY** in the window grouping to sort entries in each group.

SELECT *,
 ROW_NUMBER() OVER (ORDER BY cid)
 FROM enrolled
 ORDER BY cid

Find the student(s) with the <u>second</u> highest grade for each course.

HCMU.DB



COMMON TABLE EXPRESSIONS

Provides a way to write auxiliary statements for use in a larger query.

 \rightarrow Think of it like a temp table just for one query.

Alternative to nested queries and views.

```
WITH cteName AS (
SELECT 1
)
SELECT * FROM cteName
```

COMMON TABLE EXPRESSIONS

You can bind/alias output columns to names before the AS keyword.

```
WITH cteName (col1, col2) AS (
SELECT 1, 2
```

SELECT col1 + col2 FROM cteName

```
WITH cteName (colXXX, colXXX) AS (
SELECT 1, 2
```

SELECT * FROM cteName

CMU·DB 15-445/645 (Spring 2023) 43

COMMON TABLE EXPRESSIONS

 $\Lambda \Lambda$

Find student record with the highest id that is enrolled in at least one course.

| WI | TH cteSource (maxId) AS (SELECT MAX(sid) FROM enrolled | |
|--------------|---|--|
|) SE W | LECT name FROM student, cteSource HERE student.sid = cteSource.maxId | |
| | | |

CTE - RECURSION

Print the sequence of numbers from 1 to 10.



ECMU-DB 15-445/645 (Spring 2023

CONCLUSION

SQL is not a dead language.

You should (almost) always strive to compute your answer as a single SQL statement.

HOMEWORK #1

Write SQL queries to perform basic data analysis.

- \rightarrow Write the queries locally using SQLite.
- \rightarrow Submit them to Gradescope
- \rightarrow You can submit multiple times and use your best score.

Due: Friday, Feb 3rd @ 11:59pm

https://15445.courses.cs.cmu.edu/spring2023/homework1

SCMU·DB 15-445/645 (Spring 2023

