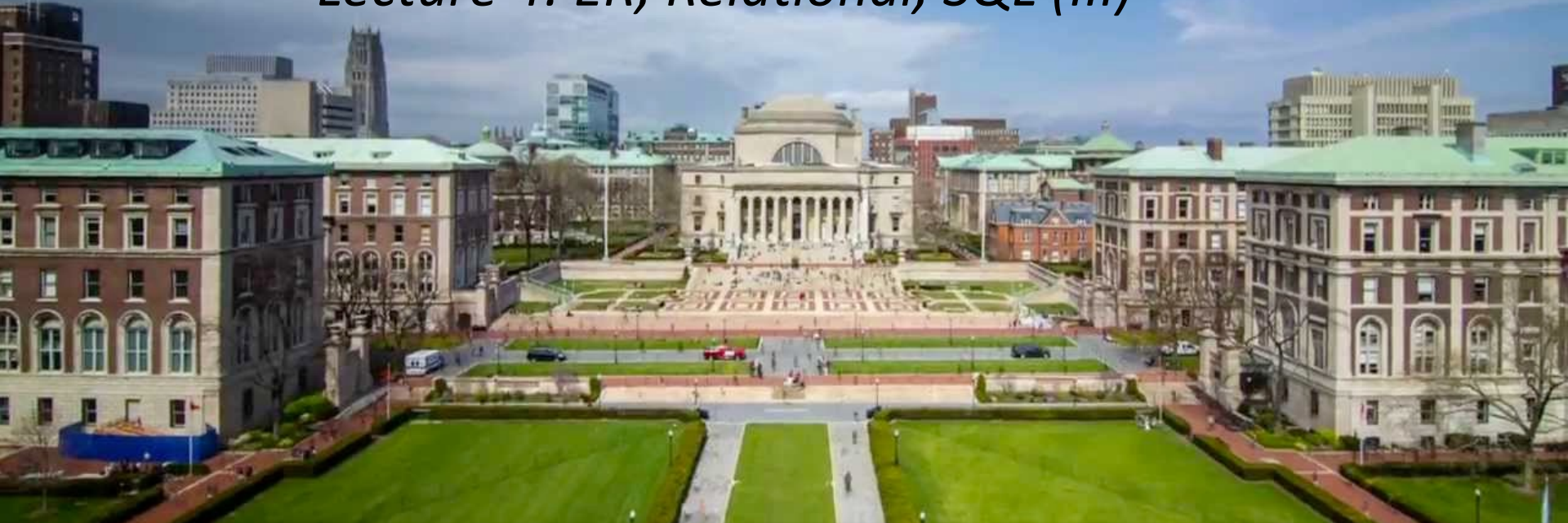


W4111 – Introduction to Databases

Section 002, Fall 2021

Lecture 4: ER, Relational, SQL (III)



Contents

Contents

- Questions, answers, discussion.
- A worked example involving:
 - ER modeling.
 - SQL DDL and data model creation.
 - Top-down and bottom-up modeling.
- A little more relational algebra and the dreaded “RelaX Calculator.”
- SQL Insert, Update and Delete.

Questions, Answers, Discussion

Systematic Treatment of NULL

From Ed Discussion:

“Confusion: When looking over the W4111_hw1_material.ipynb file, it's mentioned that "isDead is either true or NULL. NULL means unknown or not applicable". But, in the spec for the HW, it says that the column should strictly take 'Y' or 'N' values. The example output also generated confusion since it did include 'N' in the isDead column.

Question: Regardless, I understand any column can take NULL values but I just wanted to make sure we should only include 'Y' and NULL as values in our column where the latter is instead of 'N'. Should we use 'true' rather than 'Y' as the HW 1 materials file specifies as well?”

Comments:

- IMHO the question is focusing on your understanding type constraints and data cleanup. So, “Y” or “N” is fine with me.
- That aside, my solution would be *true* or NULL.
 - The data does indicate is someone is dead, but
 - The data is “old.” I cannot be sure someone is alive. They could have just died.

Codd's 12 Rules

Rule 1: Information Rule

The data stored in a database, may it be user data or metadata, must be a value of some table cell. Everything in a database must be stored in a table format.

Rule 2: Guaranteed Access Rule

Every single data element (value) is guaranteed to be accessible logically with a combination of table-name, primary-key (row value), and attribute-name (column value). No other means, such as pointers, can be used to access data.

Rule 3: Systematic Treatment of NULL Values

The NULL values in a database must be given a systematic and uniform treatment. This is a very important rule because a NULL can be interpreted as one the following – data is missing, data is not known, or data is not applicable.

Rule 4: Active Online Catalog

The structure description of the entire database must be stored in an online catalog, known as data dictionary, which can be accessed by authorized users. Users can use the same query language to access the catalog which they use to access the database itself.

Rule 5: Comprehensive Data Sub-Language Rule

A database can only be accessed using a language having linear syntax that supports data definition, data manipulation, and transaction management operations. This language can be used directly or by means of some application. If the database allows access to data without any help of this language, then it is considered as a violation.

Rule 6: View Updating Rule

All the views of a database, which can theoretically be updated, must also be updatable by the system.

Codd's 12 Rules

Rule 7: High-Level Insert, Update, and Delete Rule

A database must support high-level insertion, updation, and deletion. This must not be limited to a single row, that is, it must also support union, intersection and minus operations to yield sets of data records.

Rule 8: Physical Data Independence

The data stored in a database must be independent of the applications that access the database. Any change in the physical structure of a database must not have any impact on how the data is being accessed by external applications.

Rule 9: Logical Data Independence

The logical data in a database must be independent of its user's view (application). Any change in logical data must not affect the applications using it. For example, if two tables are merged or one is split into two different tables, there should be no impact or change on the user application. This is one of the most difficult rule to apply.

Rule 10: Integrity Independence

A database must be independent of the application that uses it. All its integrity constraints can be independently modified without the need of any change in the application. This rule makes a database independent of the front-end application and its interface.

Rule 11: Distribution Independence

The end-user must not be able to see that the data is distributed over various locations. Users should always get the impression that the data is located at one site only. This rule has been regarded as the foundation of distributed database systems.

Rule 12: Non-Subversion Rule

If a system has an interface that provides access to low-level records, then the interface must not be able to subvert the system and bypass security and integrity constraints.

Codd's 12 Rules

Rule 3: Systematic treatment of null values:

- “Null values (distinct from the empty character string or a string of blank characters and distinct from zero or any other number) are supported in fully relational DBMS for representing **missing information** and **inapplicable** information in a systematic way, independent of data type.”
- Sometimes programmers and database designers are tempted to use “special values” to indicate unknow, missing or inapplicable values.
 - String: “”, “NA”, “UNKNOWN”, ...
 - Numbers: -1, 0, -9999
- Indicators can cause confusion because you have to carefully code some SQL statements to the specific, varying choices programmers made.

NULL and Correct Answers

```
In [4]: 1 %sql describe aaaaS2lExamples.null_examples;
```

```
* mysql+pymysql://dbuser:***@localhost
3 rows affected.
```

```
Out[4]:
```

Field	Type	Null	Key	Default	Extra
name	varchar(32)	NO	PRI	None	
weight	int	YES		None	
net_worth	int	YES		None	

```
In [5]: 1 %sql select * from aaaaS2lExamples.null_examples;
```

```
* mysql+pymysql://dbuser:***@localhost
4 rows affected.
```

```
Out[5]:
```

name	weight	net_worth
Joe	100	100
Larry	0	0
Pete	None	None
Tim	200	200

Without NULL, to get a correct answer:

- I must understand the domain to determine “unknown” values or know what choice a developer made.
- Explicitly include “where weight != 0” in all statements.
- And this varies from column to column, table to table, schema to schema, etc.

```
In [7]: 1 %%sql select avg(weight) as avg_weight, avg(net_worth) as avg_net_worth
2         from aaaaS2lExamples.null_examples where name in ('Joe', 'Larry', 'Tim')

* mysql+pymysql://dbuser:***@localhost
1 rows affected.
```

```
Out[7]:
```

avg_weight	avg_net_worth
100.0000	100.0000

```
In [9]: 1 %%sql select avg(weight) as avg_weight, avg(net_worth) as avg_net_worth
2         from aaaaS2lExamples.null_examples where name in ('Joe', 'Pete', 'Tim')

* mysql+pymysql://dbuser:***@localhost
1 rows affected.
```

```
Out[9]:
```

avg_weight	avg_net_worth
150.0000	150.0000

Back to Our Excellent Question

- “When *deathYear* is null do we assume that the person is alive? What if he is dead and information is just missing?”
- Well,
 - My interpretation is that if there is no death date, *isDead* should be NULL.
 - I might expect that the person is dead if their current age would be 125, but that is just a guess.
- Also,
 - Having a column *isDead* is actually a bad design pattern in this case.
 - The value of *isDead* is *functionally dependent* on *deathYear*.
 - Database designs avoid *functional dependencies* to prevent *update anomalies*. Some could forget to update both columns when learning a fact.
 - In other scenarios, I might know that a person *isDead* but not know the *deathDate*. So, this design decision is problem/scenario specific.

Other Questions?

Worked Example:

- 1. ER Model from scratch.*
- 2. Schema definition, DDL, ... and some new concepts.*

Worked Example – Top Down

- Scenario
 - Course with complex course ID.
 - Section
 - Instructor
 - Department
 - Instructor – Department Assoc. entity with properties (role, date).
- Do ER (live) diagram in Lucidchart.
- Do schema creation
 - In DataGrip
 - Show copying statements into the notebook.



Bottom-Up: Some Sources of Information

- Course Codes:
<https://www.cc-seas.columbia.edu/sites/dsa/files/handbooks/Columbia%20Key%20to%20Course%20Listing.pdf>
- Course/Section Properties:
<http://www.columbia.edu/cu/bulletin/uwb/>
- Common Search Criteria → Indexes, Query Parameters, Data Links
<https://doc.search.columbia.edu/classes/Ferguson?instr=&name=&days=&semes=&hour=&moi=>
- Department Information:
 - Script: <https://www.columbia.edu/content/academics/departments>
 - Codes: https://academic-admin.cuit.columbia.edu/dept_code
- Specific Information:
 - Course and Instructor: <https://opendataservice.columbia.edu/api/9/json>
 - Vergil Data: <https://vergil.registrar.columbia.edu/feeds/cw.js>
 - Vergil Search: <https://vergil.registrar.columbia.edu/doc-adv-queries.php?key=ferguson&moreresults=2>

Meet in the Middle

- Show cu_info project with data and processing.
- Show notebook for loading the data.
- Show getting part of the way through parsing HTML in project

Some More Relational Algebra



Union Operation

- The union operation allows us to combine two relations
- Notation: $r \cup s$
- For $r \cup s$ to be valid.
 1. r, s must have the **same arity** (same number of attributes)
 2. The attribute domains must be **compatible** (example: 2nd column of r deals with the same type of values as does the 2nd column of s)
- Example: to find all courses taught in the Fall 2017 semester, or in the Spring 2018 semester, or in both

$$\Pi_{course_id} (\sigma_{semester="Fall" \wedge year=2017}(section)) \cup \Pi_{course_id} (\sigma_{semester="Spring" \wedge year=2018}(section))$$



Union Operation (Cont.)

- Result of:

$$\Pi_{course_id} (\sigma_{semester="Fall" \wedge year=2017}(section)) \cup$$
$$\Pi_{course_id} (\sigma_{semester="Spring" \wedge year=2018}(section))$$

<i>course_id</i>
CS-101
CS-315
CS-319
CS-347
FIN-201
HIS-351
MU-199
PHY-101

Note: The preloaded dataset on the Relax calculator is different from the most recent data referenced in the book. It is from a previous edition.



Set-Intersection Operation

- The set-intersection operation allows us to find tuples that are in both the input relations.
- Notation: $r \cap s$
- Assume:
 - r, s have the *same arity*
 - attributes of r and s are compatible
- Example: Find the set of all courses taught in both the Fall 2017 and the Spring 2018 semesters.

$$\Pi_{course_id} (\sigma_{semester="Fall" \wedge year=2017}(section)) \cap \Pi_{course_id} (\sigma_{semester="Spring" \wedge year=2018}(section))$$

- Result

<i>course_id</i>
CS-101



Set Difference Operation

- The set-difference operation allows us to find tuples that are in one relation but are not in another.
- Notation $r - s$
- Set differences must be taken between **compatible** relations.
 - r and s must have the **same** arity
 - attribute domains of r and s must be compatible
- Example: to find all courses taught in the Fall 2017 semester, but not in the Spring 2018 semester

$\Pi_{course_id} (\sigma_{semester="Fall" \wedge year=2017}(section)) -$

$\Pi_{course_id} (\sigma_{semester="Spring" \wedge year=2018}(section))$

<i>course_id</i>
CS-347
PHY-101

Same “arity”

- Same “arity”
 - Same number of columns.
 - Compatible types.
 - The i-th column in each table is from a compatible domain.
 - Student 5th column is “year.”
 - Faculty 5th column is “title”
 - Both are strings but combining them does not make sense.
- You can shape two incompatible tables using *project operations*. For example
 - π first_name, last_name, email (students)
 \cap
 π first_name, last_name, email (faculty)
 - π last_name, email, title \leftarrow 'Student' (students)
 \cup
 π last_name, email, title (faculty)

Select DB (W4111 SimpleUnion) ▼

students

id string
first_name string
last_name string
email string
year string

faculty

id string
first_name string
last_name string
email string
title string
hire_date string



The Assignment Operation

- It is convenient at times to write a relational-algebra expression by assigning parts of it to temporary relation variables.
- The assignment operation is denoted by \leftarrow and works like assignment in a programming language.
- Example: Find all instructor in the “Physics” and Music department.

$Physics \leftarrow \sigma_{dept_name=“Physics”}(instructor)$

$Music \leftarrow \sigma_{dept_name=“Music”}(instructor)$

$Physics \cup Music$

- With the assignment operation, a query can be written as a sequential program consisting of a series of assignments followed by an expression whose value is displayed as the result of the query.



The Rename Operation

- The results of relational-algebra expressions do not have a name that we can use to refer to them. The rename operator, ρ , is provided for that purpose
- The expression:

$$\rho_x(E)$$

returns the result of expression E under the name x

- Another form of the rename operation:

$$\rho_{x(A1,A2, \dots A_n)}(E)$$

Note: Assignment and rename can act a little wonky when using the calculator.



Equivalent Queries

- There is more than one way to write a query in relational algebra.
- Example: Find information about courses taught by instructors in the Physics department with salary greater than 90,000
- Query 1

$\sigma_{dept_name="Physics" \wedge salary > 90,000} (instructor)$

- Query 2

$\sigma_{dept_name="Physics"}(\sigma_{salary > 90,000} (instructor))$

- The two queries are not identical; they are, however, equivalent -- they give the same result on any database.



Equivalent Queries

- There is more than one way to write a query in relational algebra.
- Example: Find information about courses taught by instructors in the Physics department
- Query 1

$\sigma_{dept_name="Physics"}(instructor \bowtie_{instructor.ID = teaches.ID} teaches)$

- Query 2

$(\sigma_{dept_name="Physics"}(instructor)) \bowtie_{instructor.ID = teaches.ID} teaches$

- The two queries are not identical; they are, however, equivalent -- they give the same result on any database.

Insert, Update, Delete



Modification of the Database

- Deletion of tuples from a given relation.
- Insertion of new tuples into a given relation
- Updating of values in some tuples in a given relation



Deletion

- Delete all instructors

delete from *instructor*

- Delete all instructors from the Finance department

delete from *instructor*
where *dept_name* = 'Finance';

- *Delete all tuples in the instructor relation for those instructors associated with a department located in the Watson building.*

delete from *instructor*
where *dept name* in (**select** *dept name*
 from *department*
 where *building* = 'Watson');



Deletion (Cont.)

- Delete all instructors whose salary is less than the average salary of instructors

```
delete from instructor  
where salary < (select avg (salary)  
                from instructor);
```

- Problem: as we delete tuples from *instructor*, the average salary changes
- Solution used in SQL:
 1. First, compute **avg** (*salary*) and find all tuples to delete
 2. Next, delete all tuples found above (without recomputing **avg** or retesting the tuples)



Insertion

- Add a new tuple to *course*

insert into *course*

values ('CS-437', 'Database Systems', 'Comp. Sci.', 4);

- or equivalently

insert into *course* (*course_id*, *title*, *dept_name*, *credits*)

values ('CS-437', 'Database Systems', 'Comp. Sci.', 4);

- Add a new tuple to *student* with *tot_creds* set to null

insert into *student*

values ('3003', 'Green', 'Finance', *null*);



Insertion (Cont.)

- Make each student in the Music department who has earned more than 144 credit hours an instructor in the Music department with a salary of \$18,000.

```
insert into instructor
  select ID, name, dept_name, 18000
  from student
  where dept_name = 'Music' and total_cred > 144;
```

- The **select from where** statement is evaluated fully before any of its results are inserted into the relation.

Otherwise queries like

```
insert into table1 select * from table1
```

would cause problem



Updates

- Give a 5% salary raise to all instructors

```
update instructor  
  set salary = salary * 1.05
```

- Give a 5% salary raise to those instructors who earn less than 70000

```
update instructor  
  set salary = salary * 1.05  
  where salary < 70000;
```

- Give a 5% salary raise to instructors whose salary is less than average

```
update instructor  
  set salary = salary * 1.05  
  where salary < (select avg (salary)  
                  from instructor);
```



Updates (Cont.)

- Increase salaries of instructors whose salary is over \$100,000 by 3%, and all others by a 5%
 - Write two **update** statements:

```
update instructor
  set salary = salary * 1.03
  where salary > 100000;
update instructor
  set salary = salary * 1.05
  where salary <= 100000;
```
 - The order is important
 - Can be done better using the **case** statement (next slide)



Case Statement for Conditional Updates

- Same query as before but with case statement

```
update instructor  
set salary = case  
    when salary <= 100000 then salary * 1.05  
    else salary * 1.03  
end
```



Updates with Scalar Subqueries

- Recompute and update `tot_creds` value for all students

update *student S*

```
set tot_cred = (select sum(credits)  
               from takes, course  
               where takes.course_id = course.course_id and  
                   S.ID = takes.ID and  
                   takes.grade <> 'F' and  
                   takes.grade is not null);
```

- Sets `tot_creds` to null for students who have not taken any course
- Instead of **sum**(*credits*), use:

```
case  
  when sum(credits) is not null then sum(credits)  
  else 0  
end
```

Summary

- INSERT, UPDATE and DELETE are pretty straightforward.
- UPDATE and DELETE are very similar to SELECT
 - WHERE clause specifies which rows are affected.
 - The SELECT choose the columns to return.
 - The SET clause chooses and changes columns.
 - DELETE just removes the specified rows.
- INSERT, UPDATE and DELETE changes must not violate constraints, e.g.
 - INSERT a row that causes a duplicate key.
 - DELETE a referenced (target) foreign key.
 - UPDATE columns that create a duplicate key.
 - INSERT values do not include all NOT NULL columns.
- I am not going to do example now, but you have seen and will see me do examples in the context of larger examples.