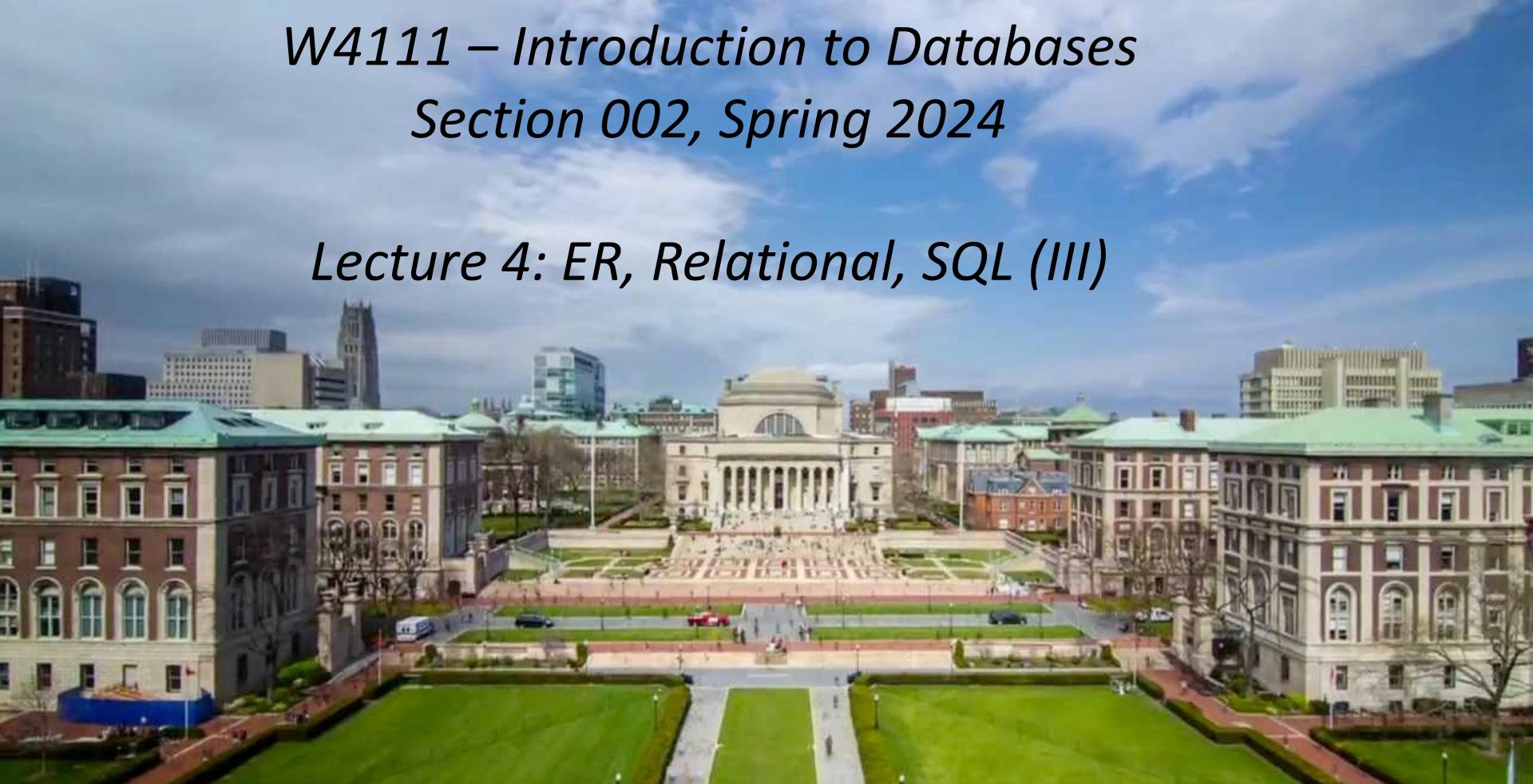


*W4111 – Introduction to Databases
Section 002, Spring 2024*

Lecture 4: ER, Relational, SQL (III)



W4111 – Introduction to Databases

Section 002, Spring 2024

Lecture 4: ER, Relational, SQL (III)

We will start in a couple of minutes.

Contents

Contents

- Some Observations
- ER (Diagram) Modeling – More Complex Scenarios: Complex Attributes
- The Relational Model and Algebra Continued (Chapter 3)
- SQL Continued (Chapter 3, 4, 5)
- Project Examples:
 - Web Applications and REST
 - Data Engineering and Visualization

Some Observations

“Create a Table” – Probable not the Best Way to Phrase but Accurate

- Like the relational algebra, SQL is closed under non-update operators, e.g.
 - SELECT
 - UNIONThe result is a table.
- There are a few types of table:
 - Base table
 - Derived table
 - Temporary table
 - View
- SELECT, UNION ... produce a derived table.

SQL 1 and 2 Tables #152



Anonymous

14 hours ago in Assignments - HW1

PIN

STAR

WATCH

63



For SQL1 and SQL2 do we actually have to create the new tables (i.e. CREATE TABLE <xxxx> AS SELECT....) or is just SELECT.... enough?

Comment Edit Delete Endorse ...

1 Answer



Joshua Zhou STAFF

13 hours ago



Just select



Comment Edit Delete Endorse ...

Add comment

- CREATE TEMPORARY TABLE creates a table that exists for the session duration (login duration).
- We will cover *views* later.

Database design, Entity-Relationship Model (Continued)

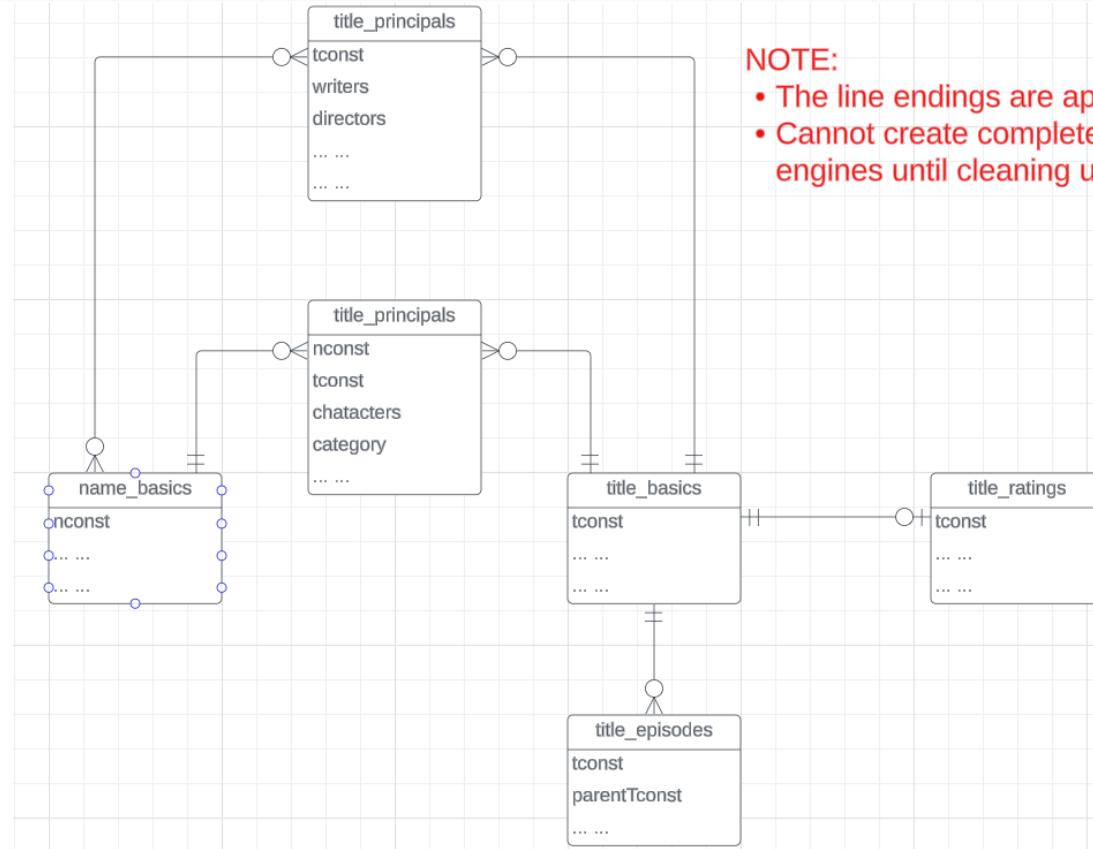
Complex Attributes

Entity Attributes

- The relational model and well-designed SQL schema have *atomic attributes*.
- There are several ways to think about attributes:
 - Simple vs Composite
 - Single Valued versus Multi-valued
 - Derived or Not Derived
- And, there can be combinations, for example a Composite, Multi-Value Attribute. Phone number is an example:
 - A phone number is a composite (+1, 914-555-1212)
 - A customer may have several: work, home, mobile,
- Examining *name_basics* from the IMDB dataset is interesting.

IMDB Free Data

- [IMDB Free Dataset](#) is a set of TSV files:
 - Title Basics
 - Title AKAS
 - Title Crew
 - Title Episodes
 - Title Principals
 - Title Ratings
 - Name Basics
- The data needs a lot of “clean up.”



Example from IMDB

Switch to notebook

- Consider name_basics

	nconst	primaryName	birthYear	deathYear	primaryProfession	knownForTitles
1	nm0000001	Fred Astaire	1899	1987	soundtrack,actor,miscellaneous	tt0050419,tt0031983,tt0072308,tt0053137
2	nm0000002	Lauren Bacall	1924	2014	actress,soundtrack	tt0071877,tt0117057,tt0037332,tt0038355
3	nm0000003	Brigitte Bardot	1934	<null>	actress,soundtrack,music_department	tt0049189,tt0056404,tt0057345,tt0054452
4	nm0000004	John Belushi	1949	1982	actor,soundtrack,writer	tt0077975,tt0072562,tt0080455,tt0078723
5	nm0000005	Ingmar Bergman	1918	2007	writer,director,actor	tt0050986,tt0060827,tt0069467,tt0050976
6	nm0000006	Ingrid Bergman	1915	1982	actress,soundtrack,producer	tt0077711,tt0038109,tt0034583,tt0036855
7	nm0000007	Humphrey Bogart	1899	1957	actor,soundtrack,producer	tt0043265,tt0034583,tt0042593,tt0037382
8	nm0000008	Marlon Brando	1924	2004	actor,soundtrack,director	tt0078788,tt0068646,tt0070849,tt0047296
9	nm0000009	Richard Burton	1925	1984	actor,soundtrack,producer	tt0061184,tt0087803,tt0057877,tt0059749
10	nm0000010	James Cagney	1899	1986	actor,soundtrack,director	tt0029870,tt0035575,tt0042041,tt0055256

- There
 - Is one composite attribute, primaryName.
 - Are two multivalued attributes: primaryProfession, knownForTitles
 - knownForTitles is also tricky, which we will see.
 - Names are also a little tricky

Generated Attributes

The Mockaroo interface shows the configuration for generating four fields:

- dept_code:** Custom List (values: COMS, EENG, ECON, MATH, DSCI)
- faculty_code:** Custom List (values: E, W, C, B, G)
- course_number:** Digit Sequence (format: ####)
- course_name:** Words (length: at least 3 but no more than 6, blank: 0%)

Buttons include "+ ADD ANOTHER FIELD" and "GENERATE FIELDS USING AI...".

Configuration parameters:
Rows: 20, Format: SQL, Table Name: MOCK_DATA, include CREATE TABLE checked.

Append Dataset: choose a dataset... ▾

Switch to Notebook

- Generated some mock data.
- Loaded into a make-believe course table in MySQL.

	A	B	C	D	E
1	dept_code	faculty_code	course_number	course_name	
2	MATH	C	3607	non ligula pellentesque ultrices phasellus id	
3	MATH	C	6017	tincidunt lacus at	
4	ECON	B	9385	velit nec nisi vulputate	
5	MATH	B	4605	congue elementum in	
6	ECON	G	2190	lectus in est risus auctor sed	
7	COMS	C	0991	dictumst morbi vestibulum	
8	ECON	C	4109	sed magna at nunc	
9	COMS	C	1845	ipsum ac tellus semper interdum mauris	
10	DSCI	B	4591	id pretium iaculis diam erat	

Relation Model and Algebra

What are all those other Symbols?

- τ order by
- γ group by
- \neg negation
- \div set division
- \bowtie natural join, theta-join
- \bowtie_l left outer join
- \bowtie_r right outer join
- \bowtie_{lf} full outer join
- \bowtie_l left semi join
- \bowtie_r right semi join
- \triangleright anti-join
- Some of these are obscure and truly bizarre
 - Division
 - Anti-Join
 - Left semi-join
 - Right semi-join
- I must look them up again every semester.
- Most SQL engines do not support them.
 - You can implement them using combinations of JOIN, SELECT, WHERE,
 - But, I cannot even remember using them in applications I have developed.
- Outer JOIN is very useful, but less common. We will cover.
- There are also some “patterns” or “terms”
 - Equijoin
 - Non-equi join
 - Natural join
 - Theta join
 -
- I may ask you to define these terms on some exams because they may be common internships/job interview questions.

Semi-Join

Semijoin (\ltimes and \bowtie) [\[edit\]](#)

The left semijoin is a joining similar to the natural join and written as $R \ltimes S$ where R and S are relations.^[b] The result is the set of all tuples in R for which there is a tuple in S that is equal on their common attribute names. The difference from a natural join is that other columns of S do not appear. For example, consider the tables *Employee* and *Dept* and their semijoin:[\[citation needed\]](#)

Employee		
Name	Empld	DeptName
Harry	3415	Finance
Sally	2241	Sales
George	3401	Finance
Harriet	2202	Production

Dept	
DeptName	Manager
Sales	Sally
Production	Harriet

Employee \ltimes Dept		
Name	Empld	DeptName
Sally	2241	Sales
Harriet	2202	Production

More formally the semantics of the semijoin can be defined as follows:

$$R \ltimes S = \{t : t \in R \wedge \exists s \in S(\text{Fun}(t \cup s))\}$$

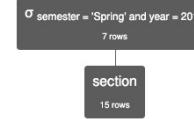
where $\text{Fun}(r)$ is as in the definition of natural join.

The semijoin can be simulated using the natural join as follows. If a_1, \dots, a_n are the attribute names of R , then

$$R \ltimes S = \Pi_{a_1, \dots, a_n} (R \bowtie S).$$

Semi-Join

- $\text{course} \ltimes (\sigma_{\text{semester}=\text{'Spring'}} \wedge \text{year}=2010) (\text{section})$
 - The common column name is *course_id*.
 - The query produces the rows from course for which there is section of the course in the Spring of 2010
- Note:
 - Walk through this in detail in RelaX.
 - We will see the corresponding SQL later.



$\sigma_{\text{semester} = \text{'Spring'}} \wedge \text{year} = 2010 (\text{section})$
Execution time: 1 ms

section.course_id	section.sec_id	section.semester	section.year	section.building	section.instructor
'CS-101'	1	'Spring'	2010	'Packard'	
'CS-315'	1	'Spring'	2010	'Watson'	
'CS-319'	1	'Spring'	2010	'Watson'	
'CS-319'	2	'Spring'	2010	'Taylor'	
'FIN-201'	1	'Spring'	2010	'Packard'	
'HIS-351'	1	'Spring'	2010	'Painter'	
'MU-199'	1	'Spring'	2010	'Packard'	

Anti-Join

Antijoin (\triangleright) [edit]

The antijoin, written as $R \triangleright S$ where R and S are [relations](#),^[c] is similar to the semijoin, but the result of an antijoin is only those tuples in R for which there is *no* tuple in S that is equal on their common attribute names.[\[citation needed\]](#)

For an example consider the tables *Employee* and *Dept* and their antijoin:

Employee		
Name	Empld	DeptName
Harry	3415	Finance
Sally	2241	Sales
George	3401	Finance
Harriet	2202	Production

Dept	
DeptName	Manager
Sales	Sally
Production	Harriet

Employee \triangleright Dept		
Name	Empld	DeptName
Harry	3415	Finance
George	3401	Finance

The antijoin is formally defined as follows:

$$R \triangleright S = \{ t : t \in R \wedge \neg \exists s \in S(Fun(t \cup s)) \}$$

or

$$R \triangleright S = \{ t : t \in R, \text{ there is no tuple } s \text{ of } S \text{ that satisfies } Fun(t \cup s) \}$$

where $Fun(t \cup s)$ is as in the definition of natural join.

The antijoin can also be defined as the [complement](#) of the semijoin, as follows:

$$R \triangleright S = R - R \ltimes S \tag{5}$$

Given this, the antijoin is sometimes called the anti-semijoin, and the antijoin operator is sometimes written as semijoin symbol with a bar above it, instead of \triangleright .

Using the sample university dataset,
Instructors that are not advisors.

- Anti-Join: $\text{instructor} \triangleright ID=i_id \text{ advisor}$
- One of many equivalent queries:
$$\pi ID, name, dept_name, salary$$

$$(\sigma i_id=NULL$$

$$(\text{instructor} \bowtie ID=i_id \text{ advisor}))$$

Set Division

- Scroll through https://en.wikipedia.org/wiki/Relational_algebra
- I am not even going to try to figure this out right now.
I already have a headache.

SQL

NULL

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 rules 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 unknown, 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 aaaaS21Examples.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 aaaaS21Examples.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  
          from aaaaS21Examples.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  
          from aaaaS21Examples.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
```



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**
 - Example: $5 + \text{null}$ returns **null**
- The predicate **is null** can be used to check for null values.
 - Example: Find all instructors whose salary is null.

```
select name  
from instructor  
where salary is null
```

- The predicate **is not null** succeeds if the value on which it is applied is not null.

Note:

- **NULL is an extremely important concept.**
- **You will find it hard to understand for a while.**



Null Values (Cont.)

- SQL treats as **unknown** the result of any comparison involving a null value (other than predicates **is null** and **is not null**).
 - Example: $5 < \text{null}$ or $\text{null} \neq \text{null}$ or $\text{null} = \text{null}$
- The predicate in a **where** clause can involve Boolean operations (**and**, **or**, **not**); thus the definitions of the Boolean operations need to be extended to deal with the value **unknown**.
 - **and** : $(\text{true and unknown}) = \text{unknown}$,
 $(\text{false and unknown}) = \text{false}$,
 $(\text{unknown and unknown}) = \text{unknown}$
 - **or**: $(\text{unknown or true}) = \text{true}$,
 $(\text{unknown or false}) = \text{unknown}$
 $(\text{unknown or unknown}) = \text{unknown}$
- Result of **where** clause predicate is treated as *false* if it evaluates to *unknown*

String Operations Between



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 substring.
 - '_ _ _' 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 "||")
 - converting from upper to lower case (and vice versa)
 - 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**

Show notebook for order by example.



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');`

Set Operations



Set Operations

- Find courses that ran in Fall 2017 or in Spring 2018

```
(select course_id from section where sem = 'Fall' and year = 2017)
union
(select course_id from section where sem = 'Spring' and year = 2018)
```

- Find courses that ran in Fall 2017 and in Spring 2018

```
(select course_id from section where sem = 'Fall' and year = 2017)
intersect
(select course_id from section where sem = 'Spring' and year = 2018)
```

- Find courses that ran in Fall 2017 but not in Spring 2018

```
(select course_id from section where sem = 'Fall' and year = 2017)
except
(select course_id from section where sem = 'Spring' and year = 2018)
```



Set Operations (Cont.)

- Set operations **union**, **intersect**, and **except**
 - Each of the above operations automatically eliminates duplicates
- To retain all duplicates use the
 - **union all**,
 - **intersect all**
 - **except all**.

Subqueries

Concepts and Examples

(Including Set Membership)



Nested Subqueries

- SQL provides a mechanism for the nesting of subqueries. A **subquery** is a **select-from-where** expression that is nested within another query.
- The nesting can be done in the following SQL query

```
select A1, A2, ..., An
  from r1, r2, ..., rm
 where P
```

as follows:

- **From clause:** r_i can be replaced by any valid subquery
- **Where clause:** P can be replaced with an expression of the form:

$B <\text{operation}> (\text{subquery})$

B is an attribute and $<\text{operation}>$ to be defined later.

- **Select clause:**

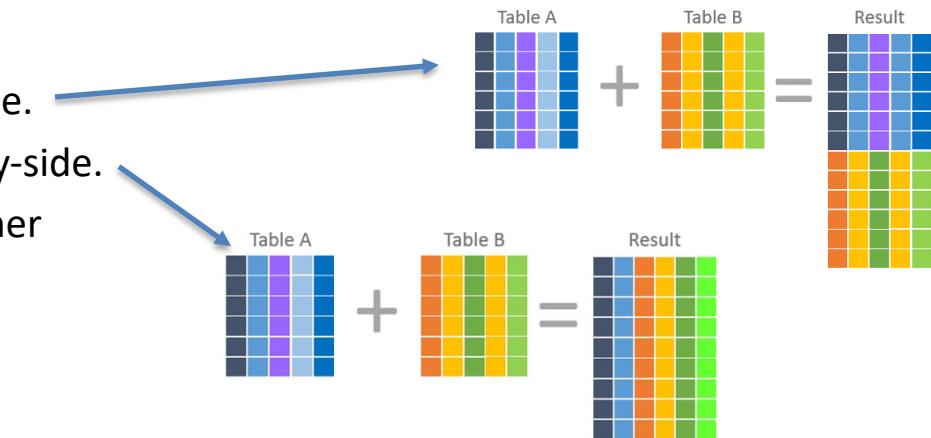
A_i can be replaced by a subquery that generates a single value.

Note:

- This is a little cryptic.
- I think I know what they mean.
- There are some operations we will see later in the material, e.g IN, EXISTS,

Nested Subquery

- The slides that come with the book have surprisingly little material on nested subqueries.
- The concept is:
 - Extremely important.
 - Students often find subqueries more confusing than joins.
 - The relationship/difference of subqueries to joins is often, initial unclear.
- We have seen:
 - Union sort of puts a table on top of a table.
 - Join puts tables sort of puts tables side-by-side.
 - Subquery enables one query to call another during execution like a subfunction.



Consider Some Tables

Takes

ID	course_id	sec_id	semester	year	grade
00128	CS-101	1	Fall	2017	A
00128	CS-347	1	Fall	2017	A-
12345	CS-101	1	Fall	2017	C
12345	CS-190	2	Spring	2017	A
12345	CS-315	1	Spring	2018	A
12345	CS-347	1	Fall	2017	A
19991	HIS-351	1	Spring	2018	B
23121	FIN-201	1	Spring	2018	C+
44553	PHY-101	1	Fall	2017	B-
45678	CS-101	1	Fall	2017	F
45678	CS-101	1	Spring	2018	B+
45678	CS-319	1	Spring	2018	B
54321	CS-101	1	Fall	2017	A-
54321	CS-190	2	Spring	2017	B+
55739	MU-199	1	Spring	2018	A-
76543	CS-101	1	Fall	2017	A
76543	CS-319	2	Spring	2018	A
76653	EE-181	1	Spring	2017	C
98765	CS-101	1	Fall	2017	C-
98765	CS-315	1	Spring	2018	B
98988	BIO-101	1	Summer	2017	A
98988	BIO-301	1	Summer	2018	None

Student

ID	name	dept_name	tot_cred
00128	Zhang	Comp. Sci.	102
12345	Shankar	Comp. Sci.	32
19991	Brandt	History	80
23121	Chavez	Finance	110
44553	Peltier	Physics	56
45678	Levy	Physics	46
54321	Williams	Comp. Sci.	54
55739	Sanchez	Music	38
70557	Snow	Physics	0
76543	Brown	Comp. Sci.	58
76653	Aoi	Elec. Eng.	60
98765	Bourikas	Elec. Eng.	98
98988	Tanaka	Biology	120

Consider a Subquery Tables

select *, (select name from student where student.id=takes.id) as name from takes;

Takes

ID	course_id	sec_id	semester	year	grade
00128	CS-101	1	Fall	2017	A
00128	CS-347	1	Fall	2017	A-
12345	CS-101	1	Fall	2017	C
12345	CS-190	2	Spring	2017	A
12345	CS-315	1	Spring	2018	A
12345	CS-347	1	Fall	2017	A
19991	HIS-351	1	Spring	2018	B
23121	FIN-201	1	Spring	2018	C+
44553	PHY-101	1	Fall	2017	B-
45678	CS-101	1	Fall	2017	F
45678	CS-101	1	Spring	2018	B+
45678	CS-319	1	Spring	2018	B
54321	CS-101	1	Fall	2017	A-
54321	CS-190	2	Spring	2017	B+
55739	MU-199	1	Spring	2018	A-
76543	CS-101	1	Fall	2017	A
76543	CS-319	2	Spring	2018	A
76653	EE-181	1	Spring	2017	C
98765	CS-101	1	Fall	2017	C-
98765	CS-315	1	Spring	2018	B
98988	BIO-101	1	Summer	2017	A
98988	BIO-301	1	Summer	2018	None

- Assume I wrote a function `find_student_name(x)`
 - Input is an `x`
 - Loops through all students and returns students with `student.ID = x`.
- The query with a subquery above is like:

`result = []`

For `t` in `takes`:

```
new_r = t + find_student_name(t.id)
result.append(new_r)
```

Switch to Notebook

Student

ID	name	dept_name	tot_cred
00128	Zhang	Comp. Sci.	102
12345	Shankar	Comp. Sci.	32
19991	Brandt	History	80
23121	Chavez	Finance	110
44553	Peltier	Physics	56
45678	Levy	Physics	46
54321	Williams	Comp. Sci.	54
55739	Sanchez	Music	38
70557	Snow	Physics	0
76543	Brown	Comp. Sci.	58
76653	Aoi	Elec. Eng.	60
98765	Bourikas	Elec. Eng.	98
98988	Tanaka	Biology	120



Nested Subqueries

- SQL provides a mechanism for the nesting of subqueries. A **subquery** is a **select-from-where** expression that is nested within another query.
- The nesting can be done in the following SQL query

```
select A1, A2, ..., An
  from r1, r2, ..., rm
 where P
```

as follows:

- **From clause:** r_i can be replaced by any valid subquery
- **Where clause:** P can be replaced with an expression of the form:

$B <\text{operation}> (\text{subquery})$

B is an attribute and $<\text{operation}>$ to be defined later.

- **Select clause:**
 A_i can be replaced by a subquery that generates a single value.

Note: Subquery MUST return

- A single scalar if in the SELECT.
- A Table if in the FROM.
- If in the WHERE:
 - Either a scalar or a table.
 - Depending on the operation.



Set Membership



Set Membership

- Find courses offered in Fall 2017 and in Spring 2018

```
select distinct course_id
from section
where semester = 'Fall' and year= 2017 and
course_id in (select course_id
from section
where semester = 'Spring' and year= 2018);
```

- Find courses offered in Fall 2017 but not in Spring 2018

```
select distinct course_id
from section
where semester = 'Fall' and year= 2017 and
course_id not in (select course_id
from section
where semester = 'Spring' and year= 2018);
```



Set Membership (Cont.)

- Name all instructors whose name is neither “Mozart” nor Einstein”

```
select distinct name  
from instructor  
where name not in ('Mozart', 'Einstein')
```

- Find the total number of (distinct) students who have taken course sections taught by the instructor with *ID* 10101

```
select count (distinct ID)  
from takes  
where (course_id, sec_id, semester, year) in  
      (select course_id, sec_id, semester, year  
from teaches  
where teaches.ID= 10101);
```

- Note: Above query can be written in a much simpler manner.
The formulation above is simply to illustrate SQL features



Set Comparison



Set Comparison – “some” Clause

- Find names of instructors with salary greater than that of some (at least one) instructor in the Biology department.

```
select distinct T.name  
from instructor as T, instructor as S  
where T.salary > S.salary and S.dept name = 'Biology';
```

- Same query using > **some** clause

```
select name  
from instructor  
where salary > some (select salary  
                      from instructor  
                      where dept name = 'Biology');
```



Definition of “some” Clause

- $F <\text{comp}> \text{some } r \Leftrightarrow \exists t \in r \text{ such that } (F <\text{comp}> t)$
Where comp can be: $<$, \leq , $>$, $=$, \neq

(5 < some

0
5
6

) = true (read: 5 < some tuple in the relation)

(5 < some

0
5

) = false

(5 = some

0
5

) = true

(5 ≠ some

0
5

) = true (since $0 \neq 5$)

$(= \text{some}) \equiv \text{in}$
However, $(\neq \text{some}) \not\equiv \text{not in}$



Set Comparison – “all” Clause

- Find the names of all instructors whose salary is greater than the salary of all instructors in the Biology department.

```
select name  
from instructor  
where salary > all (select salary  
                 from instructor  
                 where dept name = 'Biology');
```



Definition of “all” Clause

- $F <\text{comp}> \text{all } r \Leftrightarrow \forall t \in r (F <\text{comp}> t)$

(5 < all

0
5
6

) = false

(5 < all

6
10

) = true

(5 = all

4
5

) = false

(5 ≠ all

4
6

) = true (since $5 \neq 4$ and $5 \neq 6$)

$(\neq \text{all}) \equiv \text{not in}$

However, $(= \text{all}) \not\equiv \text{in}$



Test for Empty Relations

- The **exists** construct returns the value **true** if the argument subquery is nonempty.
- **exists** $r \Leftrightarrow r \neq \emptyset$
- **not exists** $r \Leftrightarrow r = \emptyset$



Use of “exists” Clause

- Yet another way of specifying the query “Find all courses taught in both the Fall 2017 semester and in the Spring 2018 semester”

```
select course_id  
from section as S  
where semester = 'Fall' and year = 2017 and  
exists (select *  
from section as T  
where semester = 'Spring' and year= 2018  
and S.course_id = T.course_id);
```

- **Correlation name** – variable S in the outer query
- **Correlated subquery** – the inner query



Use of “not exists” Clause

- Find all students who have taken all courses offered in the Biology department.

```
select distinct S.ID, S.name
from student as S
where not exists ( (select course_id
                     from course
                     where dept_name = 'Biology')
except
( select T.course_id
  from takes as T
  where S.ID = T.ID));
```

- First nested query lists all courses offered in Biology
- Second nested query lists all courses a particular student took
- Note that $X - Y = \emptyset \Leftrightarrow X \subseteq Y$
- Note: Cannot write this query using = all and its variants



Test for Absence of Duplicate Tuples

- The **unique** construct tests whether a subquery has any duplicate tuples in its result.
- The **unique** construct evaluates to “true” if a given subquery contains no duplicates .
- Find all courses that were offered at most once in 2017

```
select T.course_id  
from course as T  
where unique ( select R.course_id  
from section as R  
where T.course_id= R.course_id  
and R.year = 2017);
```



Subqueries in the From Clause



Subqueries in the Form Clause

- SQL allows a subquery expression to be used in the **from** clause
- Find the average instructors' salaries of those departments where the average salary is greater than \$42,000."

```
select dept_name, avg_salary
  from ( select dept_name, avg (salary) as avg_salary
           from instructor
          group by dept_name)
       where avg_salary > 42000;
```

- Note that we do not need to use the **having** clause
- Another way to write above query

```
select dept_name, avg_salary
  from ( select dept_name, avg (salary)
           from instructor
          group by dept_name)
       as dept_avg (dept_name, avg_salary)
    where avg_salary > 42000;
```



With Clause

- The **with** clause provides a way of defining a temporary relation whose definition is available only to the query in which the **with** clause occurs.
- Find all departments with the maximum budget

```
with max_budget (value) as
  (select max(budget)
   from department)
  select department.name
  from department, max_budget
  where department.budget = max_budget.value;
```



Complex Queries using With Clause

- Find all departments where the total salary is greater than the average of the total salary at all departments

```
with dept_total(dept_name, value) as
  (select dept_name, sum(salary)
   from instructor
   group by dept_name),
dept_total_avg(value) as
  (select avg(value)
   from dept_total)
select dept_name
from dept_total, dept_total_avg
where dept_total.value > dept_total_avg.value;
```



Scalar Subquery

- Scalar subquery is one which is used where a single value is expected
- List all departments along with the number of instructors in each department

```
select dept_name,  
       ( select count(*)  
           from instructor  
          where department.dept_name = instructor.dept_name)  
      as num_instructors  
  from department;
```

- Runtime error if subquery returns more than one result tuple

Integrity Constraints, Part 1



Integrity Constraints

- Integrity constraints guard against accidental damage to the database, by ensuring that authorized changes to the database do not result in a loss of data consistency.
 - A checking account must have a balance greater than \$10,000.00
 - A salary of a bank employee must be at least \$4.00 an hour
 - A customer must have a (non-null) phone number

DFF

- Without integrity constraints in the database, maintaining data correctness requires:
 - Lots of users know what to do and do not make mistakes.
 - Dozens of programs correctly implement constraints in the code and stay up to date on changes.
- Implementing the constraints as part of the schema eliminates many issues.



Constraints on a Single Relation

- **not null**
- **primary key**
- **unique**
- **check (P)**, where P is a predicate



Not Null Constraints

- **not null**
 - Declare *name* and *budget* to be **not null**
name varchar(20) not null
budget numeric(12,2) not null



Unique Constraints

- **unique** (A_1, A_2, \dots, A_m)
 - The unique specification states that the attributes A_1, A_2, \dots, A_m form a candidate key.
 - Candidate keys are permitted to be null (in contrast to primary keys).

Simple Example

- Consider a simple example of an entity class *major*:
 - *major(id, name, track)*
 - “id” is a uniquely generated ID
 - “name” is the major name, e.g. “Computer Science,” “Economics,”
 - “track” is a sub-track/specialty within the major, e.g. “Applications,” “AI/ML,” ...
 - “track” is optional
 - The combination of *(name, track)* is *unique*.
- **Note:** In many DBMS, this automatically creates indices for keys/constraints.
- Switch to Notebook.

```
create table if not exists majors
(
    id          int auto_increment
    primary key,
    major_name  varchar(64) not null,
    major_track varchar(64) null,
    constraint table_name_pk
        unique (major_name, major_track)
);
```



The check clause

- The **check** (P) clause specifies a predicate P that must be satisfied by every tuple in a relation.
- Example: ensure that semester is one of fall, winter, spring or summer

```
create table section
  (course_id varchar (8),
   sec_id varchar (8),
   semester varchar (6),
   year numeric (4,0),
   building varchar (15),
   room_number varchar (7),
   time_slot_id varchar (4),
   primary key (course_id, sec_id, semester, year),
   check (semester in ('Fall', 'Winter', 'Spring', 'Summer')))
```

DFF:

- We could handle the *semester check* with an *enum*.
- Switch to notebook for a slightly different example.



Referential Integrity

- Ensures that a value that appears in one relation for a given set of attributes also appears for a certain set of attributes in another relation.
 - Example: If “Biology” is a department name appearing in one of the tuples in the *instructor* relation, then there exists a tuple in the *department* relation for “Biology”.
- Let A be a set of attributes. Let R and S be two relations that contain attributes A and where A is the primary key of S. A is said to be a **foreign key** of R if for any values of A appearing in R these values also appear in S.



Referential Integrity (Cont.)

- Foreign keys can be specified as part of the SQL **create table** statement
foreign key (*dept_name*) references *department*
- By default, a foreign key references the primary-key attributes of the referenced table.
- SQL allows a list of attributes of the referenced relation to be specified explicitly.
foreign key (*dept_name*) references *department* (*dept_name*)



Cascading Actions in Referential Integrity

- When a referential-integrity constraint is violated, the normal procedure is to reject the action that caused the violation.
- An alternative, in case of delete or update is to cascade

```
create table course (
    ...
    dept_name varchar(20),
    foreign key (dept_name) references department
        on delete cascade
        on update cascade,
    ...
)
```

- Instead of cascade we can use :
 - **set null**,
 - **set default**

DFF:

- I do not like using *cascade*. I think making changes should be explicit.
- Other people disagree.
- You will get some simple practice on HW or exams.

Indexes

Concepts and Examples



Index Creation

- Many queries reference only a small proportion of the records in a table.
- It is inefficient for the system to read every record to find a record with particular value
- An **index** on an attribute of a relation is a data structure that allows the database system to find those tuples in the relation that have a specified value for that attribute efficiently, without scanning through all the tuples of the relation.
- We create an index with the **create index** command

```
create index <name> on <relation-name> (attribute);
```



Index Creation Example

- **create table student**
*(ID varchar (5),
name varchar (20) not null,
dept_name varchar (20),
tot_cred numeric (3,0) default 0,
primary key (ID))*
- **create index studentID_index on student(ID)**
- The query:

```
select *  
from student  
where ID = '12345'
```

can be executed by using the index to find the required record, without looking at all records of *student*

DFF:

- An index on (column1, column2, column3)
- Is also an index on (column1) and (column1, column2)
- But not (column2), (column3), (column2, column3).
- We will see “why” later in the semester.

Sample Projects

Projects

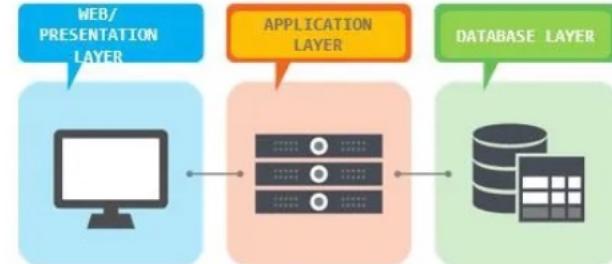
- The programming track will implement a simple, full stack web application.

Full-stack Web Developer

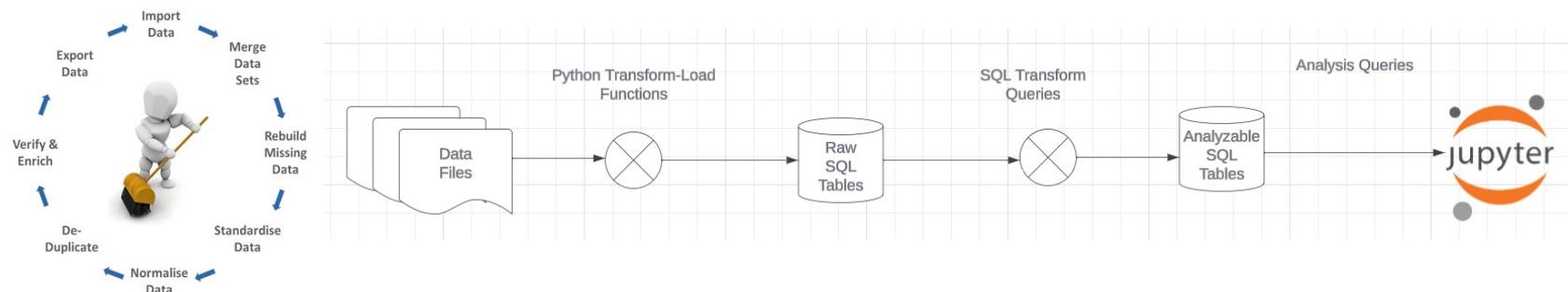
A full-stack web developer is a person who can develop both **client** and **server** software.

In addition to mastering HTML and CSS, he/she also knows how to:

- Program a **browser** (e.g. using JavaScript, jQuery, Angular, or Vue)
- Program a **server** (e.g. using PHP, ASP, Python, or Node)
- Program a **database** (e.g. using SQL, SQLite, or MongoDB)



- The non-programming track will implement a data engineering and visualization process in a Jupyter notebook.

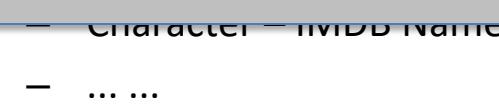


Game of Thrones

- Bottom-Up Data Mapping:
 - “Nouns” usually map to Entity/Entity Set.
 - Nouns inside other nouns often map to:
 - Attribute
 - Relationship
 - Verbs often map to relationships.
 - Adjectives usually map to properties.
- We will start with a subset of the information:
 - Game of Thrones:
 - Episodes
 - Characters
 - IMDB:
 - names_basics
 - title_basics
- Entities Sets
 - Character
 - Season
 - Episode
 - Scene
 - Location, Sublocation
 -
- Relationships
 - Character – Scene
 - Character – Character (e.g. KilledBy)
 - Season – IMDB Title
 - Character – IMDB Name
 -

Game of Thrones

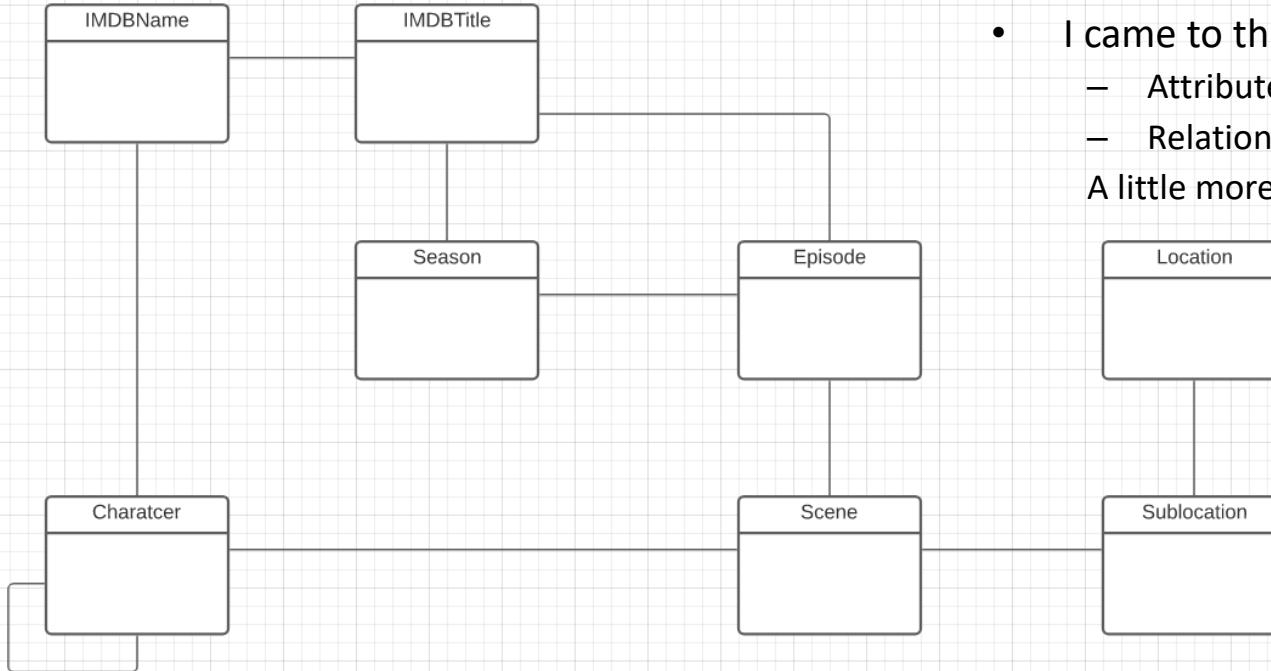
- Bottom-Up Data Mapping:
 - Entities Sets
 - We info
 - IMDB
 - 1. IMDB: <https://developer.imdb.com/non-commercial-datasets/>
 - Do not download.
 - Despite being “tiny” compared to the real world.
 - The datasets are too big for most laptops.
 - 2. Game of Thrones: <https://github.com/jeffreylancaster/game-of-thrones>
 - IMDB Data Structure
 - Character — IMDB Name
 -



Game of Thrones – Conceptual Model

Add shapes in Lucidchart ...

Add Entity Relationship Shapes



- With a little
 - Data exploration
 - Common sense
 - Judgment/experience
- I came to this conceptual model.
 - Attributes unspecified
 - Relationship required/cardinality unspecified.A little more exploration is needed.

Sample Projects

- Walkthrough of Web Apps:
 - /Users/donaldferguson/Dropbox/000/000-A-Current-Examples/W4111-FastAPI-IMDB-Solution
 - /Users/donaldferguson/Dropbox/000/000-A-Current-Examples/current-dashboard
- Data Engineering:
 - /Users/donaldferguson/Dropbox/000/000-Columbia/2024-Spring/W4111-Intro-to-Databases-Spring-2024/Lectures/s-2024-Lecture-3/More-Data-Engineering.ipynb
 - And others.