

COMS W4111: Introduction to Databases

Spring 2024, Sections 002/V02

Midterm

Introduction

This notebook contains the midterm. **Both Programming and Nonprogramming tracks should complete this.** To ensure everything runs as expected, work on this notebook in Jupyter.

- You may post **privately** on Edstem or attend OH for clarification
 - TAs will not be providing hints

Submission instructions:

- You will submit **PDF and ZIP files** for this assignment. Gradescope will have two separate assignments for these.
- For the PDF:
 - The most reliable way to save as PDF is to go to your browser's menu bar and click `File -> Print`. Switch the orientation to landscape mode, and hit save.
 - **MAKE SURE ALL YOUR WORK (CODE AND SCREENSHOTS) IS VISIBLE ON THE PDF. YOU WILL NOT GET CREDIT IF ANYTHING IS CUT OFF.** Reach out for troubleshooting.
- For the ZIP:
 - Zip a folder containing this notebook and any screenshots.
- Further submission instructions may be posted on Edstem.

Setup

```
In [ ]: %load_ext sql
        %sql mysql+pymysql://root:tzy123456@localhost
```

The sql extension is already loaded. To reload it, use:

```
%reload_ext sql
```

```
In [ ]: import pandas
        from sqlalchemy import create_engine
        engine = create_engine("mysql+pymysql://root:tzy123456@localhost")
```

Written

- You may use lecture notes, slides, and the textbook
- You may use external resources, but you must cite your sources
- As usual, keep things short

W1

Briefly explain structured data, semi-structured data, and unstructured data. Give an example of each type of data.

- structured data: follows a strict schema, defining the type and format of the data. Such as SQL databases.
- semi-structured data: doesn't reside in a rigidly defined schema but still contains tags or other markers to separate semantic elements. Such as JSON.
- unstructured data: lacks a predefined schema and the format is flexible. Such as media posts.

W2

Codd's 0th rule states:

For any system that is advertised as, or claimed to be, a relational database management system, that system must be able to manage databases entirely through its relational capabilities.

Briefly explain and give examples of how the rule applied to:

1. Metadata
2. Security

- Metadata: data that describes other data, such as the schema definitions for tables, columns, data types. For example, you can query the metadata using SQL queries against system tables or views.
- Security: encompasses identity concerns such as access control, authentication, and authorization. For example, there are tables that store user roles, permissions, and access levels to different parts of the database.

W3

Codd's 6th rule states:

All views that are theoretically updatable are also updatable by the system.

Using the following table definition, use SQL (`create view`) to define

1. Two views of the table that are not possible to update
2. One view that is possible to update

You do not need to execute the statements. We are focusing on your understanding.

```
create table student
(
    social_security_no char(9) not null primary key,
    last_name varchar(64) null,
    first_name varchar(64) null,
    enrollment_year year null,
    total_credits int null
);
```

1. Not possible to update(1): aggregated data

```
CREATE VIEW student_summary AS
SELECT enrollment_year, COUNT(*) AS total_students
FROM student
GROUP BY enrollment_year;
```

2. Not possible to update(2): complex CASE WHEN

```
CREATE VIEW student_complex AS
SELECT social_security_no,
       CASE WHEN total_credits > 100 THEN 'Senior'
            WHEN total_credits > 60 THEN 'Junior'
            WHEN total_credits > 30 THEN 'Sophomore'
            ELSE 'Freshman' END AS class_year
FROM student;
```

3. Possible to update: has simple mapping relationship

```
CREATE VIEW student_basic_info AS
SELECT social_security_no, last_name, first_name
FROM student;
```

W4

The Columbia University directory of courses uses `20241COMS4111W002` for this sections "key".

1. Is this key atomic? Explain.
 2. Explain why having non-atomic keys creates problems for indexes.
-
1. The key is not atomic. It contains different meaningful parts such as "2024" for year, "COMS" for department and "002" for section.
 2. Non-atomic keys can (1) make the index longer and take more space to store (2) cause ambiguity to extract different parts.

W5

Briefly explain the following concepts:

1. Natural join
2. Equi-join
3. Theta join
4. Left join
5. Right join
6. Outer join
7. Inner join

1. Natural join: joins two tables based on all columns with the same name and compatible data types in both tables.
2. Equi-join: merges tables based on a condition that uses equality between specified columns.
3. Theta join: allows for any conditional expression as the join condition, including `<`, `>`, `<=`, `>=`, `!=`, in addition to `=`.
4. Left join: returns all records from the left table, and the matched records from the right table. If there is no match, the result is NULL on the side of the right table.
5. Right join: returns all records from the right table, and the matched records from the left table. If there is no match, the result is NULL on the side of the left table.
6. Outer join: returns all records when there is a match in either the left or right table.
7. Inner join: returns rows when there is at least one match in both tables.

W6

The *Classic Models* database has several foreign key constraints. For instance, *orderdetails.orderNumber* references *orders.orderNumber*.

1. Briefly explain the concept of *cascading actions* relative to foreign keys.
2. How could cascading actions be helpful for the above foreign key relationship?

1. Automatically replicates the changes made to the primary key in the parent table to the corresponding foreign keys in the child table.
2. Cascading actions can simplify database maintenance and reduce errors.

W7

Give two reasons for using an associative entity to implement a relationship instead of a foreign key.

1. manage multi-to-multi relationships.
2. store relationship-specific attributes.

W8

Briefly explain how SQL is closed under its operations. Give a simple query that takes advantage of this.

The result of a SQL operation is itself a table, which can be the subject of further SQL operations.

The advantage of this is sub-query:

```
SELECT COUNT(*)
FROM (
    SELECT student_id, name
    FROM students
    WHERE grade > 90
) AS high_grade_students;
```

W9

Briefly explain the differences between:

1. Database stored procedures
 2. Database functions
 3. Database triggers
-
1. Database stored procedures: named sets of SQL instructions that perform complex operations.
 2. Database functions: user-defined routines that return a single value or a table.
 3. Database triggers: automatically executes in response to specific database events on a table or view when the requirements are met.

W10

List three benefits/use cases for defining views.

1. Simplify complex queries.
2. Restrict access to sensitive information by exposing only specific data through views.
3. Provide tailored data representations for different user needs.

Relational Algebra

- Use the [Relax calculator](#) for these questions.
- For each question, you need to show your algebra statement and a screenshot of your tree and output.
 - **For your screenshot, make sure the entire tree and output are shown.** You may need to zoom out.
- The suggestions on which relations to use are hints, not requirements.

R1

- Write a relational algebra statement that produces a relation showing **teachers that taught sections in buildings that didn't match their department's building**.
 - A section is identified by (course_id, sec_id, semester, year).
- Your output should have the following columns (names should match exactly; there should be no prefixes):
 - instructor_name
 - instructor_dept
 - course_id
 - sec_id
 - semester
 - year
 - course_building
 - dept_building
- You should use the teaches, section, instructor, and department relations.
- As an example, one row you should get is

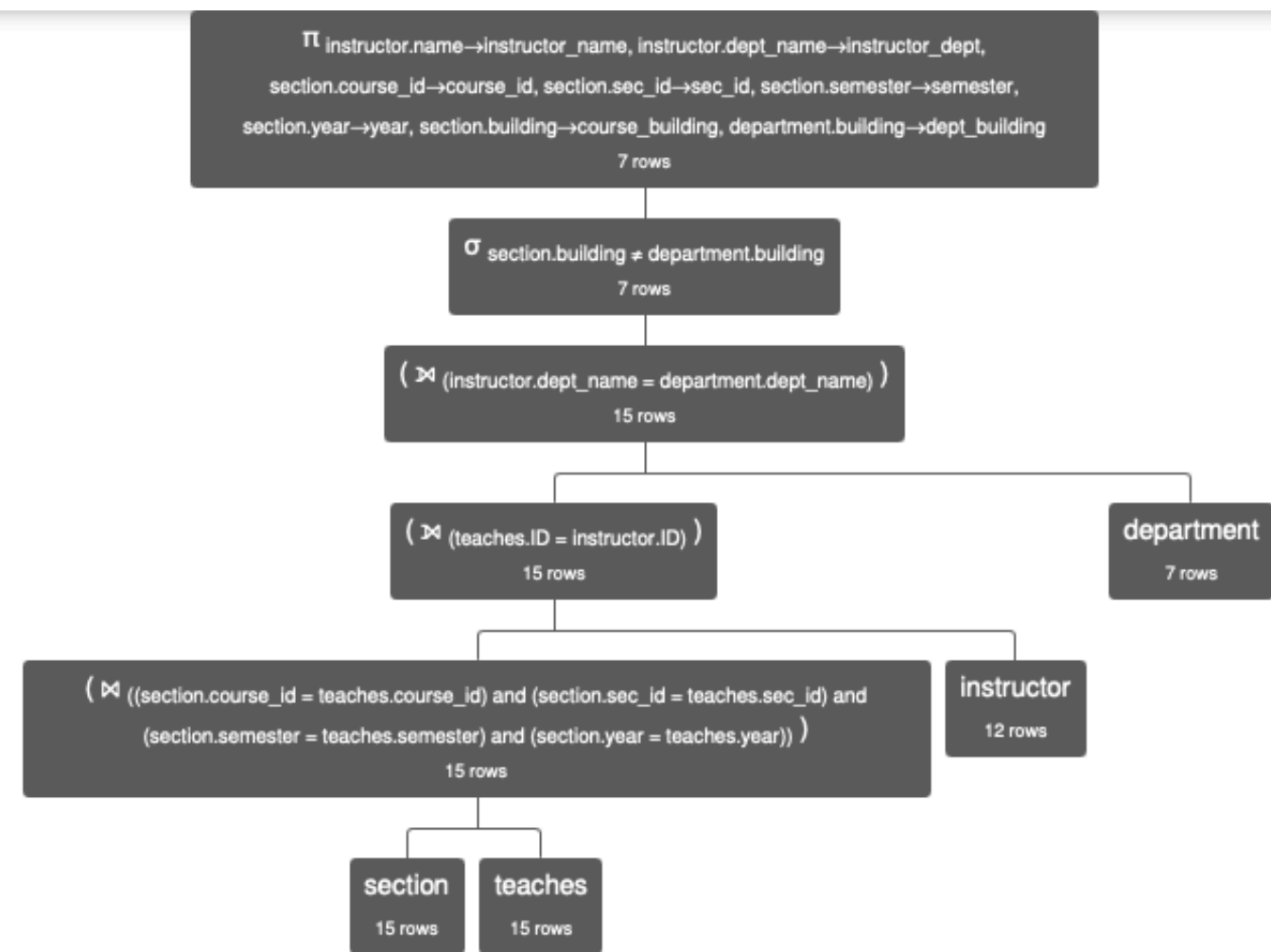
instructor_name	instructor_dept	course_id	sec_id	semester	year	course_building	dept_building
'Srinivasan'	'Comp. Sci.'	'CS-101'	1	'Fall'	2009	'Packard'	'Taylor'

- Srinivasan taught CS-101, section 1 in Fall of 2009 in the Packard building. However, Srinivasan is in the CS department, whose building is Taylor.

Algebra statement:

```
π instructor_name←instructor.name, instructor_dept←instructor.dept_name, course_id←section.course_id, sec_id←section.sec_id, semester←section.semester, year←section.year,
course_building←section.building, dept_building←department.building
(σ section.building≠department.building
(((section⋈ ((section.course_id=teaches.course_id)∧(section.sec_id=teaches.sec_id)∧
(section.semester=teaches.semester)∧(section.year=teaches.year)) teaches)
⋈ (teaches.ID=instructor.ID) instructor)
⋈ (instructor.dept_name=department.dept_name) department))
```

Execution:



π instructor.name→instructor_name, instructor.dept_name→instructor_dept, section.course_id→course_id, section.sec_id→sec_id, section.semester→semester, section.year→year, section.building→course_building, department.building→dept_building (σ section.building \neq department.building (σ section.building \neq department.building ($(($ section \bowtie ((section.course_id = teaches.course_id) and (section.sec_id = teaches.sec_id) and (section.semester = teaches.semester) and (section.year = teaches.year)) teaches) \bowtie (teaches.ID = instructor.ID) instructor) \bowtie (instructor.dept_name = department.dept_name) department)))

Execution time: 4 ms

instructor_name	instructor_dept	course_id	sec_id	semester	year	course_building	dept_building
'Crick'	'Biology'	'BIO-101'	1	'Summer'	2009	'Painter'	'Watson'
'Crick'	'Biology'	'BIO-301'	1	'Summer'	2010	'Painter'	'Watson'
'Srinivasan'	'Comp. Sci.'	'CS-101'	1	'Fall'	2009	'Packard'	'Taylor'
'Katz'	'Comp. Sci.'	'CS-101'	1	'Spring'	2010	'Packard'	'Taylor'
'Srinivasan'	'Comp. Sci.'	'CS-315'	1	'Spring'	2010	'Watson'	'Taylor'
'Katz'	'Comp. Sci.'	'CS-319'	1	'Spring'	2010	'Watson'	'Taylor'
'Wu'	'Finance'	'FIN-201'	1	'Spring'	2010	'Packard'	'Painter'

R1 Execution Result

- Some students don't have instructor advisors. Some instructors don't have student advisees.
- Write a relational algebra statement that produces a relation showing **all valid pairing between unadvised students and instructors with no advisees**.
 - A pairing is valid only if the student's department and instructor's department match.
- Your output should have the following columns (names should match exactly; there should be no prefixes):
 - instructor_name
 - student_name
 - dept_name
- You should use the `advisor`, `student`, and `instructor` relations.
- You may only use the following operators:** π , σ , $=$, \neq , \wedge (and), \vee (or), ρ , \leftarrow , \bowtie , \Join , \ltimes , \Join
 - You may not need to use all of them.
 - Notably, you may **not** use anti-join or set difference.
- As an example, one row you should get is

instructor_name	student_name	dept_name
'El Said'	'Brandt'	'History'

- El Said has no advisees, and Brandt has no advisor. They are both in the history department.
- The same instructor may show up multiple times, but the student should be different each time. Similarly, the same student may show up multiple times, but the instructor should be different each time.

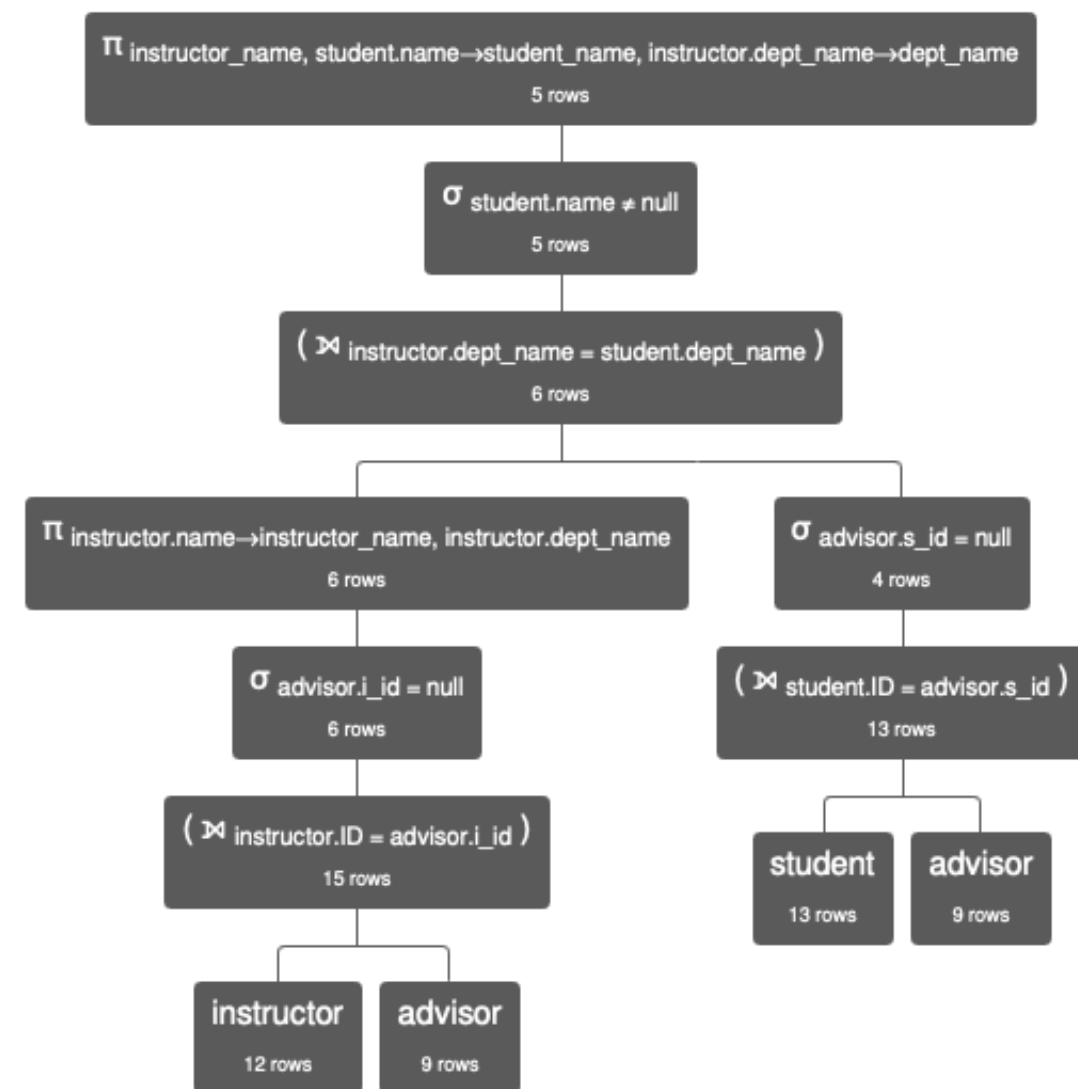
Algebra statement:

```

 $\pi$  instructor_name,student_name $\leftarrow$ student.name,dept_name $\leftarrow$ instructor.dept_name
( $\sigma$  student.name $\neq$ NULL
(( $\pi$  instructor_name $\leftarrow$ instructor.name,instructor.dept_name
( $\sigma$  advisor.i_id=NULL (instructor  $\Join$  instructor.ID=advisor.i_id advisor)))
 $\Join$  instructor.dept_name=student.dept_name
( $\sigma$  advisor.s_id=NULL (student  $\Join$  student.ID=advisor.s_id advisor))))

```

Execution:



$$\pi_{\text{instructor_name, student.name} \rightarrow \text{student_name, instructor.dept_name} \rightarrow \text{dept_name}} (\sigma_{\text{student.name} \neq \text{null}} ((\pi_{\text{instructor.name} \rightarrow \text{instructor_name, instructor.dept_name}} (\sigma_{\text{advisor.i_id} = \text{null}} (\text{instructor} \bowtie_{\text{instructor.ID} = \text{advisor.i_id}} \text{advisor}))) \bowtie_{\text{instructor.dept_name} = \text{student.dept_name}} (\sigma_{\text{advisor.s_id} = \text{null}} (\text{student} \bowtie_{\text{student.ID} = \text{advisor.s_id}} \text{advisor}))))$$

Execution time: 6 ms

instructor_name	student_name	dept_name
'Mozart'	'Sanchez'	'Music'
'El Said'	'Brandt'	'History'
'Gold'	'Snow'	'Physics'
'Califieri'	'Brandt'	'History'
'Brandt'	'Williams'	'Comp. Sci.'

R2 Execution Result

- Consider `new_section` , defined as:

`new_section = π course_id, sec_id, building, room_number, time_slot_id (section)`
- `new_section` contains sections, their time assignments, and room assignments independent of year and semester.
 - For this question, you can assume all the sections listed in `new_section` occur in the same year and semester.
 - You should copy the given definition of `new_section` to the top of your Relax calculator and treat it as a new relation.
- Write a relational algebra statement that produces a relation showing **conflicting sections**.

- Two sections conflict if they have the same `(building, room_number, time_slot_id)` .
- Your output should have the following columns (names should match exactly; there should be no prefixes):
 - `first_course_id`
 - `first_sec_id`
 - `second_course_id`
 - `second_sec_id`
 - `building`
 - `room_number`
 - `time_slot_id`
- You should use the `new_section` relation.
- Your output cannot include courses and sections that conflict with themselves, or have two rows that show the same conflict.
- Good news: I'm going to give you the correct output!

first_course_id	first_sec_id	second_course_id	second_sec_id	building	room_number	time_slot_id
'CS-190'	2	'CS-347'	1	'Taylor'	3128	'A'
'CS-319'	2	'EE-181'	1	'Taylor'	3128	'C'

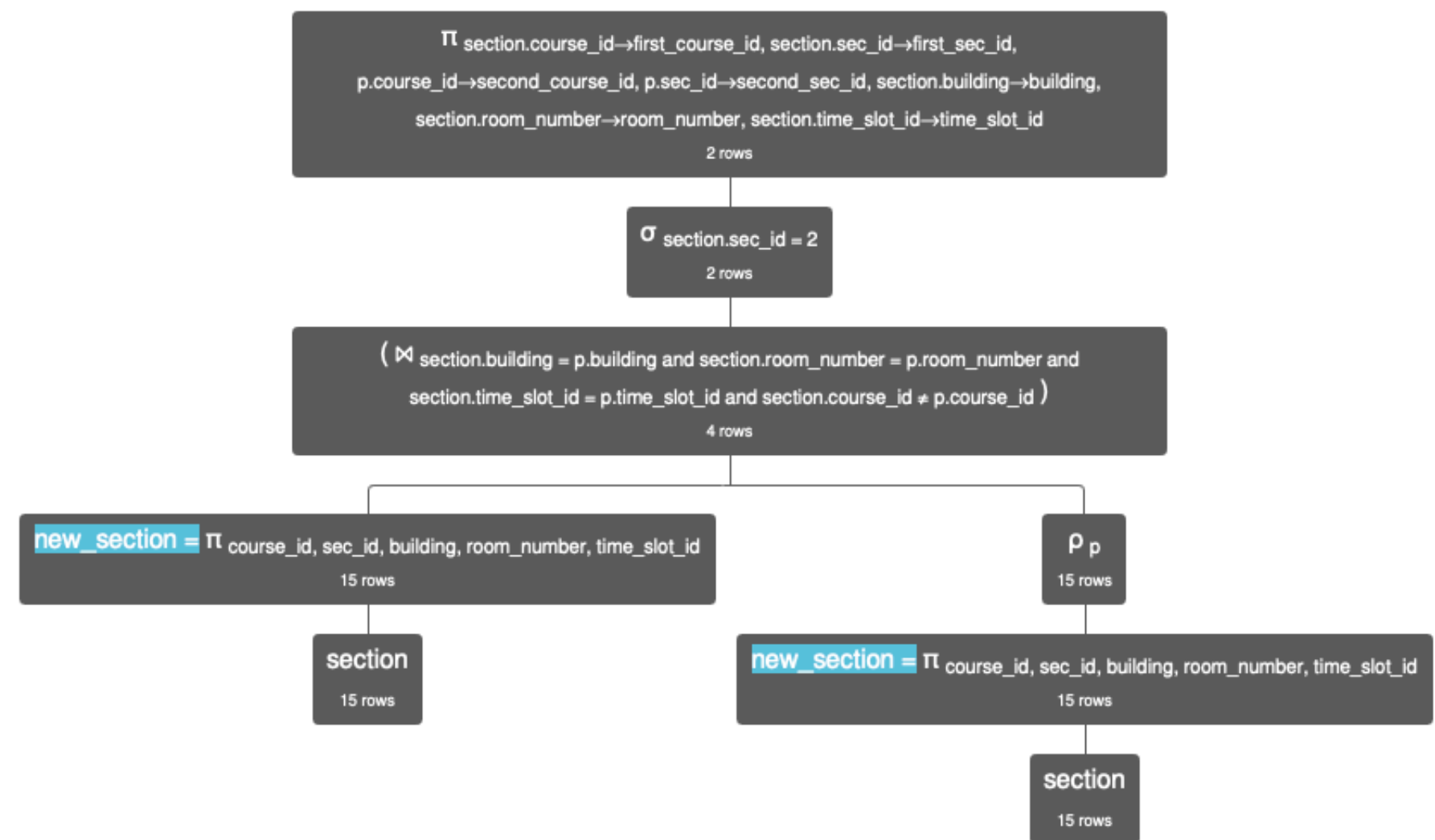
- Bad news: Your output must match mine **exactly**. The order of `first_course_id` and `second_course_id` cannot be switched.
 - Hint: You can do string comparisons in Relax using the inequality operators.

Algebra statement:

```
new_section = π course_id, sec_id, building, room_number, time_slot_id (section)

π first_course_id←section.course_id, first_sec_id←section.sec_id, second_course_id←p.course_id, second_sec_id←p.course_id, building←section.building,
room_number←section.room_number, time_slot_id←section.time_slot_id
(σ section.sec_id=2
(new_section ⋈ section.building=p.building∧section.room_number=p.room_number∧section.time_slot_id=p.time_slot_id∧section.course_id≠p.course_id p p(new_section)))
```

Execution:



Π section.course_id→first_course_id, section.sec_id→first_sec_id, p.course_id→second_course_id, p.sec_id→second_sec_id, section.building→building, section.room_number→room_number, section.time_slot_id→time_slot_id (σ section.sec_id = 2 (Π course_id, sec_id, building, room_number, time_slot_id (section) \bowtie section.building = p.building and section.room_number = p.room_number and section.time_slot_id = p.time_slot_id and section.course_id \neq p.course_id ρ p Π course_id, sec_id, building, room_number, time_slot_id (section)))

Execution time: 7 ms

first_course_id	first_sec_id	second_course_id	second_sec_id	building	room_number	time_slot_id
'CS-190'	2	'CS-347'	1	'Taylor'	3128	'A'
'CS-319'	2	'EE-181'	1	'Taylor'	3128	'C'

R3 Execution Result

ER Modeling

Definition to Model

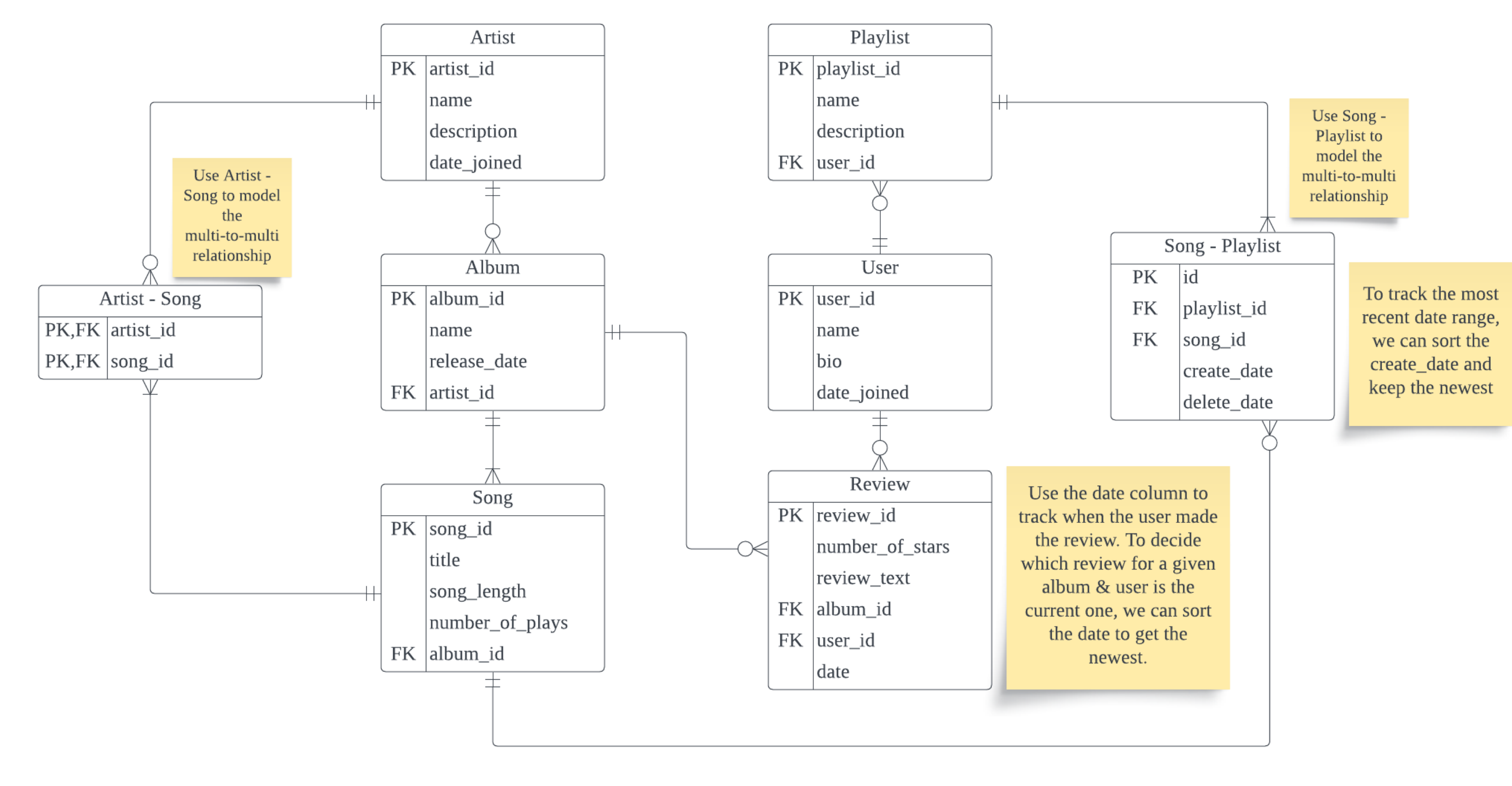
- You're in charge of creating a model for a new music app, Dotify.
- The model has the following entities:
 1. **Artist** has the properties:
 - artist_id (primary key)
 - name
 - description
 - date_joined
 2. **Album** has the properties:
 - album_id (primary key)
 - name
 - release_date
 3. **Song** has the properties:
 - song_id (primary key)
 - title
 - song_length
 - number_of_plays
 4. **User** has the properties:
 - user_id (primary key)
 - name
 - bio
 - date_joined
 5. **Review** has the properties:
 - review_id (primary key)
 - number_of_stars
 - review_text
 6. **Playlist** has the properties:
 - playlist_id (primary key)
 - name
 - description
- The model has the following relationships:
 1. **Artist–Album** : An artist can have any number of albums. An album belongs to one artist.
 2. **Album–Song** : An album can have at least one song. A song is on exactly one album.
 3. **Artist–Song** : An artist can have any number of songs. A song has at least one artist.
 4. **Album–Review** : An album can have any number of reviews. A review is associated with exactly one album.
 5. **User–Review** : A user can write any number of reviews. A review is associated with exactly one user.
 6. **User–Playlist** : A user can have any number of playlists. A playlist belongs to exactly one user.
 7. **Song–Playlist** : A song can be on any number of playlists. A playlist contains at least one song.
- Other requirements:
 1. You may **only** use the **four Crow's Foot** notations shown in class.
 2. A user can leave at most one review per album (you don't need to represent this in your diagram). However, reviews can change over time. Your model must support the ability to keep track of a user's current and previous reviews for an album as well as the dates for the reviews.
 3. Playlists can change over time. Your model must support the ability to keep track of current songs in a playlist as well as which songs were on a playlist for what date ranges.
 - You don't need to keep track of a history of when a song was on a playlist (e.g., added Jan 1, then removed Jan 2, then re-added Jan 3, then re-removed Jan 4). You can just track the most recent date range (e.g., added Jan 3, then removed Jan 4).
 4. You may not directly link many-to-many relationships. You must use an associative entity.
 5. You may (and should) add attributes to the entities and create new entities to fulfill the requirements. **Do not forget about foreign keys.**
 6. You may add notes to explain any reasonable assumptions you make, either on the Lucidchart or below.

- It would be beneficial, for instance, to document how you implemented requirements 2 and 3.

Assumptions and Documentation

- Documentations are in the diagram

Diagram:

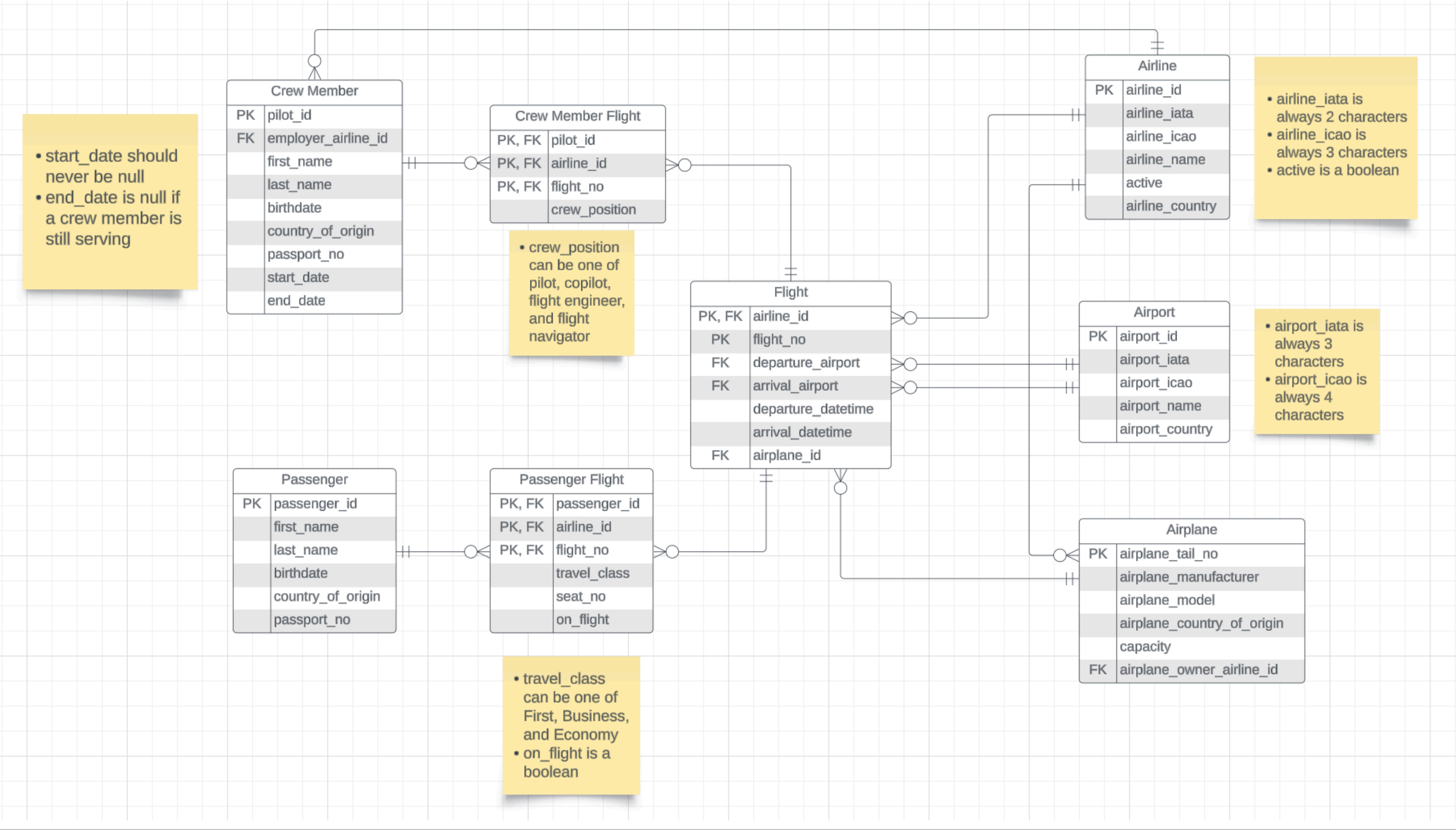


Definition to Model ER Diagram

Model to DDL

- This question tests your ability to convert an ER diagram to DDL.
- Given the ER diagram below (**not your Dotify diagram**), write `create table` statements to implement the model.
 - You should choose appropriate data types, nullness, etc.

- You are required to implement the assumptions shown in the diagram. You can document your other assumptions.
 - The required assumptions can be implemented through correct choices of data types and nullability. You aren't required to write checks or triggers for them.
- You don't need to execute your statements. You also don't need to worry about details like creating/using a database.



Model to DDL ER Diagram

Answer:

Assumptions and Documentation

- All the varchar are given a length of 255.

- `airplane_tail_no` and `airplane_id` are INT.
- Most of the attributes are NOT NULL unless specified by the notations.

```
-- Airline Table
CREATE TABLE Airline (
    airline_id INT PRIMARY KEY,
    airline_iata CHAR(2) NOT NULL,
    airline_icao CHAR(3) NOT NULL,
    airline_name VARCHAR(255) NOT NULL,
    active BOOLEAN NOT NULL,
    airline_country VARCHAR(255) NOT NULL
);

-- Airport Table
CREATE TABLE Airport (
    airport_id INT PRIMARY KEY,
    airport_iata CHAR(3) NOT NULL,
    airport_icao CHAR(4) NOT NULL,
    airport_name VARCHAR(255) NOT NULL,
    airport_country VARCHAR(255) NOT NULL
);

-- Airplane Table
CREATE TABLE Airplane (
    airplane_tail_no INT PRIMARY KEY,
    airplane_manufacturer VARCHAR(255) NOT NULL,
    airplane_model VARCHAR(255) NOT NULL,
    airplane_country_of_origin VARCHAR(255) NOT NULL,
    capacity INT NOT NULL,
    airplane_owner_airline_id INT NOT NULL,
    FOREIGN KEY (airplane_owner_airline_id) REFERENCES Airline(airline_id)
);

-- Flight Table
CREATE TABLE Flight (
    airline_id INT NOT NULL,
    flight_no INT NOT NULL,
    departure_airport INT NOT NULL,
    arrival_airport INT NOT NULL,
    departure_datetime DATETIME NOT NULL,
    arrival_datetime DATETIME NOT NULL,
    airplane_id INT NOT NULL,
    PRIMARY KEY (airline_id, flight_no),
    FOREIGN KEY (airline_id) REFERENCES Airline(airline_id),
    FOREIGN KEY (departure_airport) REFERENCES Airport(airport_id),
    FOREIGN KEY (arrival_airport) REFERENCES Airport(airport_id),
    FOREIGN KEY (airplane_id) REFERENCES Airplane(airplane_tail_no)
);

-- Crew Member Table
CREATE TABLE CrewMember (
    pilot_id INT PRIMARY KEY,
    employer_airline_id INT NOT NULL,
    first_name VARCHAR(255) NOT NULL,
    last_name VARCHAR(255) NOT NULL,
    birthdate DATE NOT NULL,
    country_of_origin VARCHAR(255) NOT NULL,
    passport_no VARCHAR(255) NOT NULL,
    start_date DATE NOT NULL,
    end_date DATE,      -- Can be NULL if the crew member is still serving
);
```

```
        FOREIGN KEY (employer_airline_id) REFERENCES Airline(airline_id)
    );

-- Passenger Table
CREATE TABLE Passenger (
    passenger_id INT PRIMARY KEY,
    first_name VARCHAR(255) NOT NULL,
    last_name VARCHAR(255) NOT NULL,
    birthdate DATE NOT NULL,
    country_of_origin VARCHAR(255) NOT NULL,
    passport_no VARCHAR(255) NOT NULL
);

-- Crew Member Flight Table
CREATE TABLE CrewMemberFlight (
    pilot_id INT,
    airline_id INT,
    flight_no INT,
    crew_position ENUM('pilot', 'copilot', 'flight engineer', 'flight navigator'),
    PRIMARY KEY (pilot_id, airline_id, flight_no),
    FOREIGN KEY (pilot_id) REFERENCES CrewMember(pilot_id),
    FOREIGN KEY (airline_id, flight_no) REFERENCES Flight(airline_id, flight_no)
);

-- Passenger Flight Table
CREATE TABLE PassengerFlight (
    passenger_id INT,
    airline_id INT,
    flight_no INT,
    travel_class ENUM('First', 'Business', 'Economy'),
    seat_no INT NOT NULL,
    on_flight BOOLEAN NOT NULL,
    PRIMARY KEY (passenger_id, airline_id, flight_no),
    FOREIGN KEY (passenger_id) REFERENCES Passenger(passenger_id),
    FOREIGN KEY (airline_id, flight_no) REFERENCES Flight(airline_id, flight_no)
);
```

Data and Schema Cleanup

Setup

- There are several issues with the `classicmodels` schema. Two issues are:
 - Having programs or users enter country names for `customers.country` is prone to error.
 - `products.productCode` is clearly not an atomic value.
- The following code does the following:
 1. Creates a schema for this question
 2. Creates copies of `classicmodels.customers` and `classicmodels.products`
 3. Loads a table of [ISO country codes](#)

In []: `%%sql`

```
drop schema if exists classicmodels_midterm;
create schema classicmodels_midterm;
use classicmodels_midterm;

create table customers as select * from classicmodels.customers;
create table products as select * from classicmodels.products;
```

```
* mysql+pymysql://root:***@localhost
4 rows affected.
1 rows affected.
0 rows affected.
122 rows affected.
110 rows affected.
```

Out[]: []

```
In [ ]: iso_df = pandas.read_csv('./wikipedia-iso-country-codes.csv')
iso_df.to_sql('countries', schema='classicmodels_midterm',
             con=engine, index=False, if_exists="replace")
```

Out[]: 246

In []: %%sql

```
alter table countries
  change `English short name` lower case `short_name` varchar(64) null;

alter table countries
  change `Alpha-2 code` alpha_2_code char(2) null;

alter table countries
  change `Alpha-3 code` alpha_3_code char(3) not null;

alter table countries
  change `Numeric code` numeric_code smallint unsigned null;

alter table countries
  change `ISO 3166-2` iso_text char(13) null;

alter table countries
  add primary key (alpha_3_code);

select * from countries limit 10;
```

```
* mysql+pymysql://root:***@localhost
246 rows affected.
246 rows affected.
246 rows affected.
246 rows affected.
246 rows affected.
0 rows affected.
10 rows affected.
```


Out[]:

short_name	alpha_2_code	alpha_3_code	numeric_code	iso_text
Aruba	AW	ABW	533	ISO 3166-2:AW
Afghanistan	AF	AFG	4	ISO 3166-2:AF
Angola	AO	AGO	24	ISO 3166-2:AO
Anguilla	AI	AIA	660	ISO 3166-2:AI
Åland Islands	AX	ALA	248	ISO 3166-2:AX
Albania	AL	ALB	8	ISO 3166-2:AL
Andorra	AD	AND	20	ISO 3166-2:AD
Netherlands Antilles	AN	ANT	530	ISO 3166-2:AN
United Arab Emirates	AE	ARE	784	ISO 3166-2:AE
Argentina	AR	ARG	32	ISO 3166-2:AR

DE1

- There are four values in `customers.country` that do not appear in `countries.short_name`.
- Write a query that finds these four countries.
 - Hint: Norway should be one of these countries.

In []:

```
%%sql

select distinct
  cs.country as country
from customers cs
  left join countries ct
    on cs.country = ct.short_name
where
  ct.short_name IS NULL
```

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

Out[]:

country
USA
Norway
UK
Russia

DE2

- `Norway` actually does appear in `countries.short_name`. The reason it appeared in DE1 is because there are two spaces after the name (`Norway__` instead of `Norway`).
- The mapping for the other countries is:

customers.country	countries.short_name
USA	United States
UK	United Kingdom
Russia	Russian Federation

- Write `update table` statements to correct the values in `customers.country` so that all the values in that attribute appear in `countries.short_name`.

```
In [ ]: %%sql

update customers
set country = "United States"
where country = "USA";

update customers
set country = "United Kingdom"
where country = "UK";

update customers
set country = "Russian Federation"
where country = "Russia";

update customers
set country = "Norway"
where country = "Norway ";

* mysql+pymysql://root:***@localhost
36 rows affected.
5 rows affected.
1 rows affected.
2 rows affected.
```

Out[]: []

DE3

- The PK of `countries` is `alpha_3_code`. We want that as a FK in `customers`.
1. Create a column `customers.iso_code`
 2. Set `customers.iso_code` as a FK that references `countries.alpha_3_code`
 3. Fill `customers.iso_code` with the appropriate data based on `customers.country`
 4. Drop `customers.country`
 5. Create a view `customers_country` of form `(customerNumber, customerName, country, iso_code)`

Bonus point: I would ask you to create an index on `customers.iso_code`, but this is actually already done for us. When was an index created on `customers.iso_code`?

Answer

```
In [ ]: %%sql

ALTER TABLE customers
ADD COLUMN iso_code VARCHAR(3);

ALTER TABLE customers
ADD CONSTRAINT fk_customers_countries
FOREIGN KEY (iso_code) REFERENCES countries(alpha_3_code);

UPDATE customers
SET iso_code = (
    SELECT alpha_3_code
    FROM countries
    WHERE countries.short_name = customers.country
);

ALTER TABLE customers
DROP COLUMN country;
```

```
CREATE VIEW customers_country AS
SELECT c.customerNumber, c.customerName, ct.short_name AS country, c.iso_code
FROM customers c
JOIN countries ct ON c.iso_code = ct.alpha_3_code;
```

```
* mysql+pymysql://root:***@localhost
0 rows affected.
122 rows affected.
122 rows affected.
0 rows affected.
0 rows affected.
```

Out[]: []

DE4

- To test your code, output a table that shows the number of customers from each country.
- You should use your `customers_country` view.
- Your table should have the following attributes:
 - `country_iso`
 - `number_of_customers`
- Order your table from greatest to least `number_of_customers`.
- Show only the first 10 rows.

```
In [ ]: %%sql

select
    iso_code as country_iso,
    count(*) as number_of_customers
from
    customers_country
group by iso_code
order by number_of_customers desc
limit 10;
```

```
* mysql+pymysql://root:***@localhost
10 rows affected.
```

Out[]:

country_iso	number_of_customers
USA	36
DEU	13
FRA	12
ESP	7
GBR	5
AUS	5
ITA	4
NZL	4
FIN	3
CAN	3

DE5

- `products.productCode` appears to be 3 separate values joined by an underscore.
 - I have no idea what the values mean, but let's pretend we do know for the sake of this question.
- Write `alter table` statements to create 3 new columns: `product_code_letter`, `product_code_scale`, and `product_code_number`.
 - Choose appropriate data types. `product_code_letter` should always be a single letter.

In []: `%%sql`

```
alter table products
add column product_code_letter char(1),
add column product_code_scale varchar(16),
add column product_code_number varchar(16);
```

```
* mysql+pymysql://root:***@localhost
0 rows affected.
```

Out[]: []

DE6

- As an example, for the product code `S18_3856`, the product code letter is `S`, the product code scale is `18`, and the product code number is `3856`.
 - I know the product code scale doesn't always match `products.productScale`. Let's ignore this for now.
1. Populate `product_code_letter`, `product_code_scale`, and `product_code_number` with the appropriate values based on `productCode`.
 2. Set the PK of `products` to `(product_code_letter, product_code_scale, product_code_number)`.
 3. Drop `productCode`.

In []: `%%sql`

```
update products
set product_code_letter = (
    select SUBSTR(productCode,1,1)
)
where product_code_letter IS NULL;

update products
set product_code_scale = (
    select SUBSTR(SUBSTRING_INDEX(productCode,'_',1),2,10)
)
where product_code_scale IS NULL;

update products
set product_code_number = (
    select SUBSTRING_INDEX(productCode,'_',-1)
)
where product_code_number IS NULL;

alter table products
add constraint pk
primary key (product_code_letter, product_code_scale, product_code_number);

alter table products
drop productCode;
```

```
* mysql+pymysql://root:***@localhost
110 rows affected.
110 rows affected.
110 rows affected.
0 rows affected.
0 rows affected.
```

Out[]: []

DE7

- To test your code, output a table that shows the products whose `product_code_scale` doesn't match `productScale` .
- Your table should have the following attributes:
 - `product_code_letter`
 - `product_code_scale`
 - `product_code_number`
 - `productScale`
 - `productName`
- Order your table on `productName` .

```
In [ ]: %%sql

select
    product_code_letter,
    product_code_scale,
    product_code_number,
    productScale,
    productName
from products
where product_code_scale != SUBSTR(productScale,3,10)
order by productName;
```

* mysql+pymysql://root:***@localhost
6 rows affected.

Out []:

product_code_letter	product_code_scale	product_code_number	productScale	productName
S	24	3856	1:18	1956 Porsche 356A Coupe
S	24	4620	1:18	1961 Chevrolet Impala
S	12	3148	1:18	1969 Corvair Monza
S	700	2824	1:18	1982 Camaro Z28
S	700	3167	1:72	F/A 18 Hornet 1/72
S	18	2581	1:72	P-51-D Mustang

SQL

- Use the `classicmodels` database for these questions.
- The suggestions on which tables to use are hints, not requirements.
- All your answers should be a single select statement. **You may not create a new table.**
 - Subqueries (selects within a select) and the `with` keyword are fine. Just don't use the `create` keyword.

```
In [ ]: %%sql use classicmodels
```

* mysql+pymysql://root:***@localhost
0 rows affected.

Out []: []

SQL1

- Write a query that produces a table of form `(productName, productLine, productVendor, totalRevenue)` .

- Attribute names should match exactly.
 - The `totalRevenue` for a product is the sum of `quantityOrdered*priceEach` across all the rows the product appears in in `orderdetails` .
 - You should consider all orders, regardless of `orders.status` .
- Only include products with `totalRevenue` greater than \$150,000.

- Order your output on `totalRevenue` descending.

- You should use the `products` and `orderdetails` tables.

In []: %%sql

```
with one as (  
    select  
        productCode,  
        sum(quantityOrdered * priceEach) as totalRevenue  
    from orderdetails  
    group by productCode  
)  
  
select  
    p.productName as productName,  
    p.productLine as productLine,  
    p.productVendor as productVendor,  
    o.totalRevenue as totalRevenue  
from one o  
left join products p  
on o.productCode = p.productCode  
where totalRevenue > 150000  
order by totalRevenue desc;
```

* mysql+pymysql://root:***@localhost
6 rows affected.

Out[]:

	productName	productLine	productVendor	totalRevenue
	1992 Ferrari 360 Spider red	Classic Cars	Unimax Art Galleries	276839.98
	2001 Ferrari Enzo	Classic Cars	Second Gear Diecast	190755.86
	1952 Alpine Renault 1300	Classic Cars	Classic Metal Creations	190017.96
	2003 Harley-Davidson Eagle Drag Bike	Motorcycles	Red Start Diecast	170686.00
	1968 Ford Mustang	Classic Cars	Autoart Studio Design	161531.48
	1969 Ford Falcon	Classic Cars	Second Gear Diecast	152543.02

SQL2

- Write a query that produces a table of form `(productCode, productName, productVendor, customerCount)`.
 - Attribute names should match exactly.
 - `customerCount` is the number of **distinct** customers that have bought the product.
 - Note that the same customer may buy a product multiple times. This only counts as one customer in the product's `customerCount` .
 - You should consider all orders, regardless of `status` .
- Order your table from largest to smallest `customerCount` , then on `productCode` alphabetically.

- Only show the first 10 rows.

- You should use the `orders` and `orderdetails` tables.

```
In [ ]: %%sql

select
  p.productCode as productCode,
  p.productName as productName,
  p.productVendor as productVendor,
  tmp.customerCount as customerCount
from products p
left join(
  select
    productCode,
    count(distinct customerNumber) as customerCount
  from orderdetails od
  left join orders o
    on od.orderNumber = o.orderNumber
  group by productCode
) tmp
on p.productCode = tmp.productCode
order by customerCount desc, productCode
limit 10
;
```

* mysql+pymysql://root:***@localhost
10 rows affected.

Out []:

productCode	productName	productVendor	customerCount
S18_3232	1992 Ferrari 360 Spider red	Unimax Art Galleries	40
S10_1949	1952 Alpine Renault 1300	Classic Metal Creations	27
S10_4757	1972 Alfa Romeo GTA	Motor City Art Classics	27
S18_2957	1934 Ford V8 Coupe	Min Lin Diecast	27
S72_1253	Boeing X-32A JSF	Motor City Art Classics	27
S10_1678	1969 Harley Davidson Ultimate Chopper	Min Lin Diecast	26
S10_2016	1996 Moto Guzzi 1100i	Highway 66 Mini Classics	26
S18_1662	1980s Black Hawk Helicopter	Red Start Diecast	26
S18_1984	1995 Honda Civic	Min Lin Diecast	26
S18_2949	1913 Ford Model T Speedster	Carousel DieCast Legends	26

SQL3

- Write a query that produces a table of form (customerName, month, year, monthlyExpenditure, creditLimit) .
 - Attribute names should match exactly.
 - monthlyExpenditure is the total amount of payments made by a customer in a specific month and year based on the payments table.
 - Some customers have never made any payments ever. For these customers, monthlyExpenditure should be 0. month and year can be null.
- Only show rows where monthlyExpenditure exceeds creditLimit or the customer has never made any payments ever (so month and year should be null for these rows).
- Order your table on monthlyExpenditure descending, then on customerName alphabetically.
- Only show the first 10 rows.
- You should use the payments and customers tables.

```
In [ ]: %%sql
```

```
with one as (  
    select customerNumber,  
           YEAR(paymentDate) as year,  
           MONTH(paymentDate) as month,  
           sum(amount) as monthlyExpenditure_old  
    from payments  
    group by customerNumber, year, month  
)  
  
select  
    customerName,  
    month,  
    year,  
    if(monthlyExpenditure_old is NULL, 0, monthlyExpenditure_old) as monthlyExpenditure,  
    creditLimit  
from customers c  
left join one o  
    on c.customerNumber = o.customerNumber  
where if(monthlyExpenditure_old is NULL, 0, monthlyExpenditure_old) > c.creditLimit OR if(monthlyExpenditure_old is NULL, 0, monthlyExpenditure_old) = 0  
order by monthlyExpenditure desc, customerName  
limit 10  
;
```

* mysql+pymysql://root:***@localhost
10 rows affected.

Out[]:

	customerName	month	year	monthlyExpenditure	creditLimit
	Dragon Souveniers, Ltd.	12	2003	105743.00	103800.00
	American Souvenirs Inc	None	None	0.00	0.00
	ANG Resellers	None	None	0.00	0.00
	Anton Designs, Ltd.	None	None	0.00	0.00
	Asian Shopping Network, Co	None	None	0.00	0.00
	Asian Treasures, Inc.	None	None	0.00	0.00
	BG&E Collectables	None	None	0.00	0.00
	Cramer Spezialitäten, Ltd	None	None	0.00	0.00
	Der Hund Imports	None	None	0.00	0.00
	Feuer Online Stores, Inc	None	None	0.00	0.00

SQL4

- Write a query that produces a table of form (productCode, productName, productLine, productVendor, productDescription) .
 - Attribute names should match exactly.
- **You should only keep products that have never been ordered by a French customer.**
 - You should consider all orders, regardless of status .
- Order your table on productCode .
- You should use the customers , orders , and orderdetails tables.

In []: %%sql

```
select  
    productCode,
```



```
    productName,
    productLine,
    productVendor,
    productDescription
from products
where productCode not in
(select distinct productCode
from orders o
left join customers c
    on o.customerNumber = c.customerNumber
left join orderdetails od
    on o.orderNumber = od.orderNumber
where c.country = 'France')
order by productCode
;
```

* mysql+pymysql://root:***@localhost
2 rows affected.

Out []:

productCode	productName	productLine	productVendor	productDescription
S18_3233	1985 Toyota Supra	Classic Cars	Highway 66 Mini Classics	This model features soft rubber tires, working steering, rubber mud guards, authentic Ford logos, detailed undercarriage, opening doors and hood, removable split rear gate, full size spare mounted in bed, detailed interior with opening glove box
S18_4027	1970 Triumph Spitfire	Classic Cars	Min Lin Diecast	Features include opening and closing doors. Color: White.

SQL5

- A customer can have a sales rep employee.
- Corporate is deciding which employees to give raises to.
 - A raise is given for the reason customers if an employee has 8 or more customers.
 - A raise is given for the reason orders if the total number of orders made by customers associated with an employee is 30 or greater.
 - You should consider all orders, regardless of status .
 - A raise is given for the reason both if both conditions above are true.
- Write a query that produces a table of form (firstName, lastName, totalCustomers, totalCustomerOrders, raiseBecause) .
 - Attribute names should match exactly.
 - firstName and lastName are for the employee.
 - totalCustomers is the total number of customers associated with an employee.
 - totalCustomerOrders is the total number of orders made by customers associated with an employee.
 - raiseBecause is one of customers , orders , and both .
- Your table should only show employees eligible for raises, i.e., raiseBecause should not be null.
- Order your table on firstName .
- You should use the customers , orders , and employees tables.

In []:

```
%%sql

with one as (
    select
        c.customerNumber as customerNumber,
        count(*) as totalOrders,
        salesRepEmployeeNumber
    from customers c
    left join orders o
        on c.customerNumber = o.customerNumber
```

```

    group by c.customerNumber
), two as (
    select
        firstName,
        lastName,
        count(*) as totalCustomers,
        sum(totalOrders) as totalCustomerOrders
    from employees e
    left join one
        on e.employeeNumber = one.salesRepEmployeeNumber
    group by e.employeeNumber
)

select
    firstName,
    lastName,
    totalCustomers,
    totalCustomerOrders,
    case
        when totalCustomerOrders >= 30 and totalCustomers >= 8 then 'both'
        when totalCustomerOrders >= 30 then 'orders'
        when totalCustomers >= 8 then 'customers'
        else 'wrong'
    end as raiseBecause
from two
where totalCustomerOrders >= 30 or totalCustomers >= 8
order by firstName
;
```

* mysql+pymysql://root:***@localhost
6 rows affected.

Out[]:

firstName	lastName	totalCustomers	totalCustomerOrders	raiseBecause
Barry	Jones	9	25	customers
George	Vanauf	8	22	customers
Gerard	Hernandez	7	43	orders
Larry	Bott	8	22	customers
Leslie	Jennings	6	34	orders
Pamela	Castillo	10	31	both