

# COMS W4111-002, Fall 23: Take Home Midterm

---

## Overview

## Instructions

**Due Date: Sunday, October 22, 2023 at 11:59pm**

You have one week to complete the take home midterm. All of the work must be your own; you may not work in groups or teams. You may use outside sources so long as you cite them and provide links. A citation would be of the form

I used this source for this question: <https://stackoverflow.com/questions/298739/what-is-the-difference-between-a-schema-and-a-table-and-a-database>.

Points will be taken off for any answers that are verbose. Try to stay between 2-3 sentences for definitions and 5 sentences for longer questions.

There is a [pinned thread](#) on Edstem for corrections and calrifications. **Students are responsible for monitoring the thread. Do not waste our time by asking a question that has an answer.**

You may post **privately** on Edstem or attend OH for clarification questions. TAs will not be providing hints or help, except to clarify questions.

Students have a bad habit of posting questions, requests for help, etc., on forums and Edstem. We often see questions like,

- "This is my SQL for question 2, does anyone know why it does not work?"
- "I almost get the answer but it is off a little. Any suggestions?"

**Do not do things like the above.** That is a serious academic violation.

## Submission Instructions

The TAs will post [submission instructions](#) on EdSTEM.

## Environment Setup

You may need to change the MySQL userID and password in some of the cells below to match your configuration.

```
In [2]: %load_ext sql
```

```
In [3]: %sql mysql+pymysql://root:dbuserdbuser@localhost
```

```
In [4]: %sql SELECT 1
```

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

```
Out[4]: 1  
1
```

```
In [5]: import pandas as pd  
from sqlalchemy import create_engine
```

```
In [6]: sql_engine = create_engine("mysql+pymysql://root:dbuserdbuser@localhost")
```

```
In [15]: from IPython.display import Image
```

## Written Questions

**Note:** You may have to look in lecture notes, slides, slides associated with the recommended textbook, or online research to answer these questions. You will have to do online search and show initiatives in your careers to answer questions more complex than the ones in this exam.

You should cite online references/links. You do not need to cite the lectures, slides from lectures, or textbook material.

### W1

Provide a short (two or three sentence) definition/description of the following terms. Some of these concepts do not have a single, precise, agreed definition. You may find slight differences in your research. Focus on the concept and grading will be flexible.

1. Super Key
2. Candidate Key
3. Primary Key
4. Alternate Key
5. Unique Key
6. Natural Key
7. Surrogate Key

8. Substitute Key

9. Foreign Key

10. External Key

### Answer

- Super Key: A set of one or more attributes that can be used to uniquely identify a record in a table.
- Candidate Key: A minimal super key that can be used to uniquely identify a record.
- Primary Key: A arbitrarily chosen candidate key, usually immutable and contains no null values.
- Alternate Key: A candidate key that is not chosen as the primary key.
- Unique Key: An attribute or set of attributes that can uniquely identify a record within a table, similar to a primary key, but it allows for one null value.
- Natural Key: A type of unique key that has intrinsic meaning in the real world.
- Surrogate Key: An artificial, system-generated unique value used to uniquely identify a record in a table, having no semantic meaning.
- Substitute Key: Another term for a surrogate key.
- Foreign Key: An attribute or set of attributes in one table that refers to the primary key in another table, establishing a relationship between the two tables.
- External Key: A key linking a table to external system or dataset.

I used this source for this question: <https://blog.devart.com/surrogate-key-in-sql.html#:~:text=Conclusion-,What%20is%20a%20surrogate%20key%20in%20SQL,object%20generates%20this%20key%20itself.>

## W2

1. Define the concept of *immutable* column.
2. Why do some sources recommend that a primary key should be composed of immutable columns?

### Answer

- An "immutable column" in a table means that once a value is inserted or initially set for that column, it should not be modified thereafter.
- Reasons: Changing primary keys can lead to confusion or errors in record referencing. It can disrupt foreign key relationships in other tables. Changing primary keys can disrupt consistent record identification, therefore disrupt data accuracy and enabling efficient data retrieval.

## W3

Views are a powerful concept in relational database management systems. List and briefly explain 3 benefits of/reasons for creating a view.

### Answer

- To hide original table information such as sensitive columns
- Since not everyone good at complex SQL, view make it easy to query
- When changes are made to physical level, views won't change and make sure upper level won't break

## W4

Briefly explain the concepts of *procedural* language and *declarative* language. SQL is a declarative language. What are some advantages of a declarative language over a procedural language?

### Answer

Procedural Language: A programming approach where the user specifies a sequence of steps to achieve a desired result. Declarative Language: A programming approach where the user defines what they want to achieve without explicitly outlining the steps to get there.

Pros:

- Simple, easier to read and understand
- Shorter and more concise, reducing the chances of errors and making code maintenance easier

Cons:

- Limited control over the execution process, hard to fine-tune performance or specify the order of operations in detail
- Debugging can be more challenging because we can't inspect each step of execution directly

## W5

The following diagram is a simple representation of the architecture of a Jupyter notebook using MySQL. Is this a two-tier architecture or a three-tier architecture? Explain your answer briefly.



### Answer

The diagram represents a three-tier architecture.

Presentation Tier: This is represented by the "User" and "Browser". It's where the user interacts with the system.

Logic Tier: Comprising the "Notebook server" and "Kernel", this is where processing occurs and business logic resides.

Database Tier: Represented by "MySQL", this tier handles data storage and retrieval.

## W6

What is the difference between the database schema and the database instance? Do you use DDL for schema or instances? Do you use DML for schemas or instances?

### Answer

Schema: A blueprint that defines the structure, organization, and constraints of a database.

Instance: An instance represents the actual data stored in a database at a specific point in time.

DDL is used for schema, DML is used for instances.

## W7

The lecture slides and a previous homework defined a convention/notation for documenting the schema for a relation in the relational model. Use the notation to define the schema for a relation with the following columns:

- product\_category
- product\_code
- product\_name
- description

The primary key is composed of product\_category and product\_code .

The following cell shows how to format text. Double click on the cell to see the source.

*Something how\_to\_underline, some other stuff* (1)

Answer

*Product(product\_category, product\_code, product\_name, description)* (2)

## W8

Briefly define and explain:

- Natural join
- Equi-join
- Theta join
- Self-join

Answer

Natural Join: A join that returns rows when there is a match in both tables based on all columns with the same name and data types. It automatically matches columns and discards duplicate columns.

Equi-join: A join on the equality between two columns, typically from two different tables. It uses the '=' operator to match rows based on the specified columns.

Theta Join: A more general form of join that uses any kind of binary comparison operator, to match rows based on a condition between two columns, like =, <, >, etc.

Self-join: A join operation in which a table is joined with itself. It's useful when comparing rows within the same table or finding relationships within the same dataset.

## W9

Briefly explain the difference between a *unique (key) constraint* and a *primary key constraint*?

Answer

- A table can have only one primary key constraint, but multiple unique constraints.
- The columns defined as the primary key cannot contain NULL values, while unique constraint columns can have NULLs.

## W10

Briefly explain the difference between a *column (data) type* and the *column's domain*.

Answer

Domain: Set of all possible values for an attribute or column, defines the permissible range of values.

Data Type: Specifies the data format and associated constraints.

Domain is a constraint on value of data, while data type is a constraint on the format of data.



# Entity Relationship Model

- This question tests transforming a high-level description of a data model into a more concrete *logical ER* diagram. You will produce a logical ER diagram using Lucidchart. You should use Crow's Foot notation and conventions we have used in lectures and examples.
- The data model is a simple representation of a university.
- The model has the following entity types:
  - School:
    - School code, e.g., "SEAS," "GSAS," "LAW," ... ..
    - School name, e.g., "School of Engineering and Applied Science."
  - Department:
    - Department code, e.g., "COMS," "MATH," "ECON," ... ..
    - Department name, e.g., "Department of Computer Science."
  - Faculty:
    - UNI
    - last name
    - first name
    - email
    - title, e.g., "Professor," "Adjunct Professor," ... ..
  - Student:
    - UNI
    - last name
    - first name
    - email
  - Course:
    - Course number is a composite key, e.g., "COMSW4111" is
      - Dept. code "COMS"

- Faculty code "W"
  - Course number "4111"
  - Course title
  - Course description
  - Section:
    - Call number
    - Course number
    - Year
    - Semester
    - Section
- A Faculty has complex states and relationships.
  - A Faculty can have a role relative to a Department. The possible roles are `chair`, `professor`, `adjunct` and `emeritus`.
  - Roles change over time. The data model must support the ability to handle current roles and previous roles, and the dates for the roles.
- A Student has a relationship to Section.
  - The possible roles are `(Enrolled, Waitlist, Dropped, TA)`
  - The student has a current role, and there can be only one current role.
  - The data model must support the ability to handle roles changing over time and retain information about prior roles.
- A Faculty *may* teach a Section. All sections have exactly one Faculty.
- The relationship between Department and School is many-to-many. Each Department is in at least one school and each school has at least one department.

**Notes:**

1. There is no single correct answer to this question. You will have to make some design decisions and assumptions. You should document your decisions and assumptions. You can do this by putting a note/comment in your diagram.

2. You do not have to worry about **isA** relationships.
3. You do not have to document or worry about attribute types.
4. The ER diagram must be implementable in the relational/SQL model.

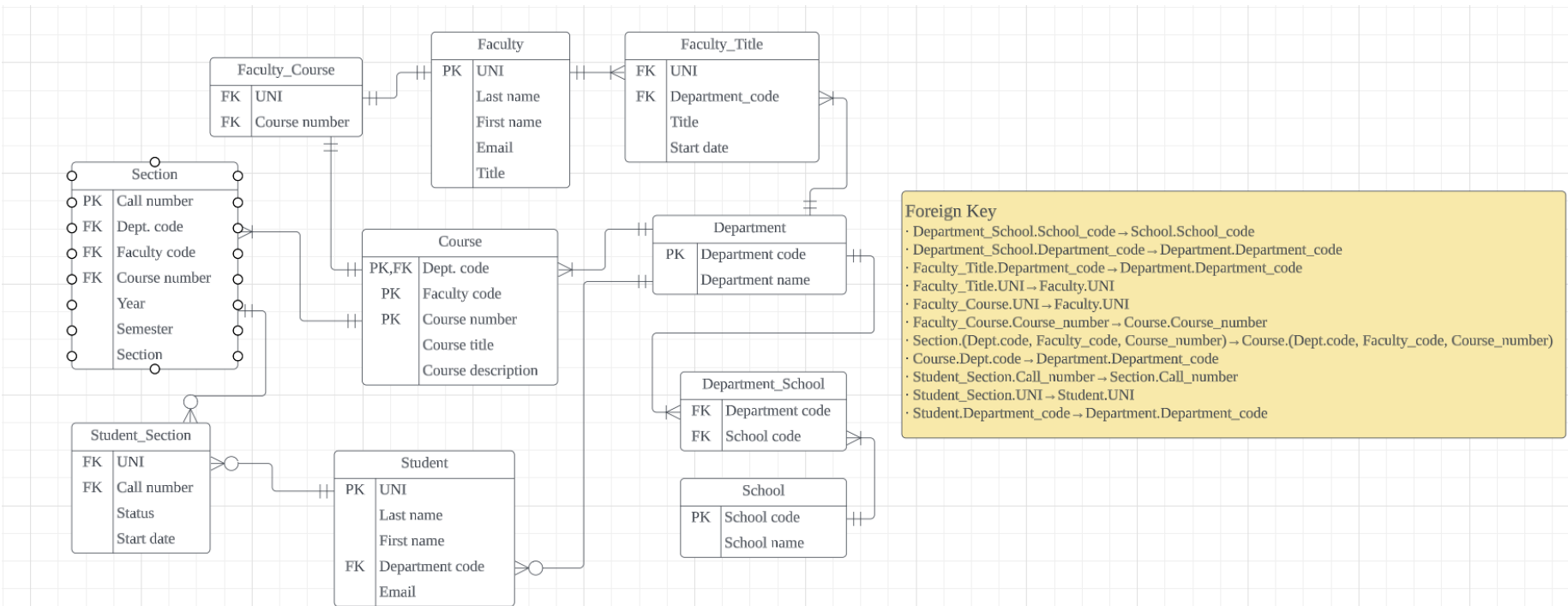
### Answer

#### Assumptions:

1. Each student only belongs to one department
2. Each course is only under one department
3. A department can have zero to many students
4. A student can register zero to many sections(courses)

In [17]: Image('ER.png')

Out[17]:



## Relational Algebra

# R1

Use the [RelaX Calculator and the Silberschatz - UniversityDB](#) for this question.

Two time slots  $X$  and  $Y$  obviously overlap if:

1. They are not the same time slot, i.e. do not have the same *time\_slot\_id*.
2. They have at least one lecture on *the same day*, the start hour for  $X$  is before the start hour for  $Y$ , and the end hour for  $X$  is after the start hour for  $Y$ .
3. To make the question easier, you do not need to consider minutes in computing overlap but must show minutes in the result.

Write the relational algebra expression that identifies obviously overlapping time slots, and only lists overlapping pairs of time slots once.

Your output must match the answer below.



## Answer

```

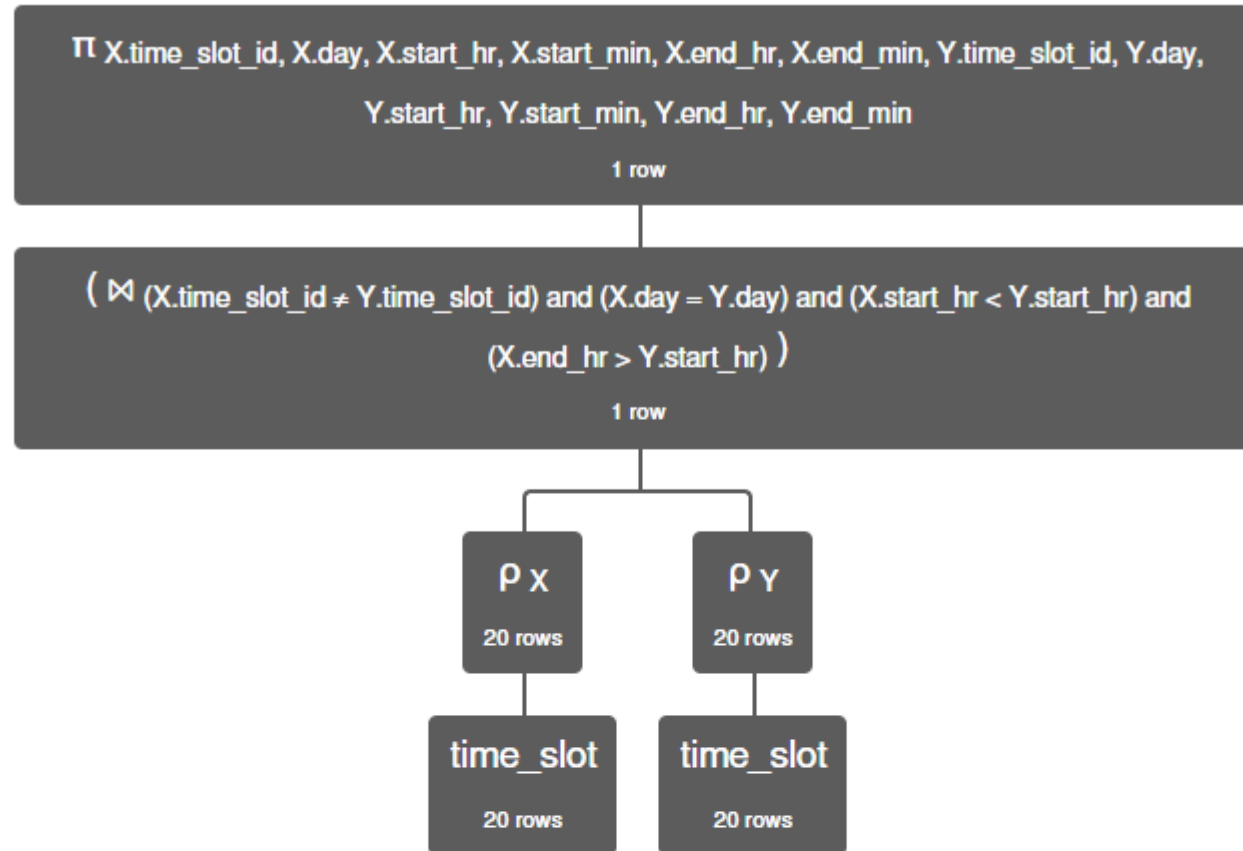
_____  $\pi$  X.time_slot_id,
_____ X.day,
_____ X.start_hr,
_____ X.start_min,
_____ X.end_hr,
_____ X.end_min,
_____ Y.time_slot_id,
_____ Y.day,
_____ Y.start_hr,
_____ Y.start_min,
_____ Y.end_hr,
_____ Y.end_min(
_____  $\rho$  X(time_slot)  $\bowtie$ 
_____ (X.time_slot_id  $\neq$  Y.time_slot_id)  $\wedge$ 
_____ (X.day = Y.day)  $\wedge$ 
_____ (X.start_hr < Y.start_hr)  $\wedge$ 

```

```
_____(X.end_hr > Y.start_hr).  
____p Y(time_slot))
```

In [18]: `Image('R1 result.png')`

Out[18]:



$\Pi$  X.time\_slot\_id, X.day, X.start\_hr, X.start\_min, X.end\_hr, X.end\_min, Y.time\_slot\_id, Y.day, Y.start\_hr, Y.start\_min, Y.end\_hr, Y.end\_min (  $\rho_X$  ( time\_slot )  $\bowtie$  (X.time\_slot\_id  $\neq$  Y.time\_slot\_id) and (X.day = Y.day) and (X.start\_hr < Y.start\_hr) and (X.end\_hr > Y.start\_hr)  $\rho_Y$  ( time\_slot ) )

Execution time: 0 ms

X.time_slot_id	X.day	X.start_hr	X.start_min	X.end_hr	X.end_min	Y.time_slot_id	Y.day	Y.start_hr	Y.start_min	Y.end_hr	Y.end_min
----------------	-------	------------	-------------	----------	-----------	----------------	-------	------------	-------------	----------	-----------

---

'H'	'W'	10	0	12	30	'C'	'W'	1
-----	-----	----	---	----	----	-----	-----	---

## R2

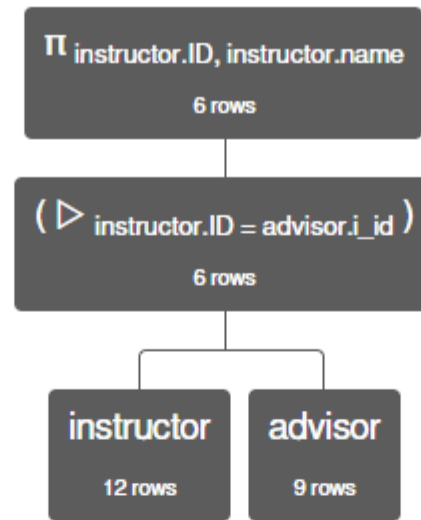
1. You **may not** use the subtraction operator `-` to write this query.
2. Produce a relation that:
  - Has column names `instructor_ID`, `instructor_name`.
  - Contains the `ID` and `name` of instructors who do not advise any students.

### Answer

$\pi$  instructor.ID, instructor.name (instructor  $\triangleright$  instructor.ID = advisor.i\_id advisor)

In [19]: `Image('R2_result.png')`

Out[19]:



$$\pi \text{ instructor.ID, instructor.name } ( \text{instructor} \triangleright \text{instructor.ID} = \text{advisor.i\_id} \text{ advisor} )$$

Execution time: 10 ms

instructor.ID	instructor.name
12121	'Wu'
15151	'Mozart'
32343	'El Said'
33456	'Gold'
58583	'Califieri'



# SQL Schema and DDL

## Objective

- You have a logical datamodel ER diagram (see below).
- You need to use DDL to define a schema that realizes the model. Name your schema `f23_w4111_midterm_medical`.
- Logical models are not specific enough for direct implementation. This means that:
  - You will have to assign concrete types to columns, and choose things like `GENERATED`, `DEFAULT`, etc.
  - You may have to decompose a table into two tables, or extract common attributes from multiple tables into a single, referenced table.
  - Implementing the relationships may require adding columns and foreign keys, associative entities, etc.
  - You may have to make other design and implementation choices. **This means that there is no single correct answer.**
- You may need to look do some self-study to find information about concepts, e.g. `GENERATED`. Cite your sources.

## ER Diagram



### ER Diagram

Answer

Design Decisions, Notes, etc.

I used this source for this question: <https://www.simplilearn.com/tutorials/sql-tutorial/auto-increment-in-sql#:~:text=The%20auto%20increment%20in%20SQL,for%20every%20record%20you%20add.>

Here is my logic in building this schema:

- All the common columns between table are not specified in the ER diagram, therefore should be added specifically. If A to B is one to many relationship and B to A is one to one relationship, refer table B to table A when adding foreign keys by including PK of table A to table B and establish the FK between two columns.
- Since one insurance company can insure many patients, insured person column should not appear in Insurance Company table, but should add insurance company id as a column in Patient table.
- Auto increment are used to generate receipt no in Payment since no column in other table is referring this column.

### DDL

- Execute your DDL in the cell below. You may use DataGrip or other tools to help build the DDL statements.
- You can copy and paste the SQL `CREATE/ALTER TABLE` statements below, but you MUST execute the statements in the notebook.

In [126...

```
%%sql
# DDL in cells below.

DROP SCHEMA IF EXISTS f23 w4111 midterm medical;
CREATE SCHEMA f23 w4111 midterm medical;
USE f23 w4111 midterm medical;

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

Out[126]:

```
[]
```

In [127...

```
%%sql

create table if not exists f23 w4111 midterm medical.insurance company
(
    company id varchar(50) not null
```

```

        primary key,
        phone no    varchar(12)    null,
        street1     varchar(255)   null,
        city        varchar(255)   null,
        state       varchar(50)    null,
        postal code  varchar(50)    null
    );

create table if not exists f23 w4111 midterm medical.patient
(
    patient ID      varchar(50)    not null
        primary key,
    last name       varchar(9)     not null,
    first name      varchar(9)     not null,
    street1         varchar(255)   null,
    city            varchar(255)   null,
    state           varchar(50)    null,
    postal code     varchar(10)    null,
    insur company id varchar(50)    null,
    constraint patient insurance company company id fk
        foreign key (insur company id) references f23 w4111 midterm medical.insurance company (company id)
);

create table if not exists f23 w4111 midterm medical.physician
(
    physician ID    varchar(50)          not null
        primary key,
    last name       varchar(9)           not null,
    first name      varchar(9)           not null,
    physician type enum ('NP', 'MD', 'DO') not null
);

create table if not exists f23 w4111 midterm medical.appointment
(
    appt ID        varchar(50)    not null
        primary key,
    appt date      date           not null,
    appt time      time           not null,
    appt duration decimal(10, 2) not null,
    appt reason    text           null,
    physician ID   varchar(50)    not null,
    patient ID     varchar(50)    not null,
    constraint appointment patient patient ID fk
        foreign key (patient ID) references f23 w4111 midterm medical.patient (patient ID),

```

```

    constraint appointment_physician_physician ID fk
    foreign key (physician ID) references f23 w4111 midterm medical.physician (physician ID)
);

create table if not exists f23 w4111 midterm medical.bill
(
    bill no          varchar(50)          not null
    primary key,
    amount insured   decimal(10, 2)       not null,
    amount not insured decimal(10, 2)     not null,
    bill total       decimal(10, 2)       null,
    bill date        date                 null,
    bill status      varchar(9) default 'Pending' null,
    appt ID          varchar(50)          null,
    constraint bill_appointment_appt ID fk
    foreign key (appt ID) references f23 w4111 midterm medical.appointment (appt ID)
);

create table if not exists f23 w4111 midterm medical.payment
(
    receipt no       int auto increment
    primary key,
    paid amount      decimal(10, 2)       not null,
    paid date        date                 null,
    paid type        enum ('CC', 'Check', 'Transfer') not null,
    bill no          varchar(9)           null,
    insur company id varchar(50)          null,
    patient ID       varchar(50)          null,
    constraint payment_bill_bill no fk
    foreign key (bill no) references f23 w4111 midterm medical.bill (bill no),
    constraint payment_insurance_company_company id fk
    foreign key (insur company id) references f23 w4111 midterm medical.insurance company (company id),
    constraint payment_patient_patient ID fk
    foreign key (patient ID) references f23 w4111 midterm medical.patient (patient ID)
);

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

```

Out[127]: []

In [ ]: \_

## Complex SQL

### Birth Countries and Death Countries

- **Note:** Use the instructions from HW1 Part 2 Nonprogramming to load the Lahman's Baseball Database data into MySQL if you have not already done so.
- In Lahman's Baseball Database table `people`, there is information about people's `birthCountry` and `deathCountry`.
- There are countries in which at least one person was born but in which no person has died.
- Write a query that produces a table of the form:
  - `birthCountry`
  - `no_of_births`, which is the total number of births in the country.
- The table contains all rows for countries in which there with births but no deaths.
- You may create indexes to improve performance.

### Answer

In [96]: %%sql

USE lahmans hw1;

```
WITH death AS (  
    SELECT deathCountry, COUNT(*) AS no of deaths  
    FROM people  
    GROUP BY deathCountry  
)  
  
birth AS (  
    SELECT birthCountry, COUNT(*) AS no of births  
    FROM people  
    GROUP BY birthCountry  
)  
  
SELECT birthCountry, no of births  
FROM birth  
LEFT JOIN death ON birth.birthCountry = death.deathCountry  
WHERE death.deathCountry IS NULL AND birthCountry != 'None';  
  
* mysql+pymysql://root:***@localhost  
0 rows affected.  
35 rows affected.
```

Out[96]: **birthCountry** **no of births**

<u>Russia</u>	<u>9</u>
<u>Curacao</u>	<u>16</u>
<u>Colombia</u>	<u>28</u>
<u>Nicaragua</u>	<u>15</u>
<u>Germany</u>	<u>45</u>
<u>Norway</u>	<u>3</u>
<u>Italy</u>	<u>7</u>
<u>South Korea</u>	<u>27</u>
<u>Czech Republic</u>	<u>7</u>
<u>Aruba</u>	<u>6</u>
<u>Sweden</u>	<u>4</u>
<u>Hong Kong</u>	<u>1</u>
<u>Afghanistan</u>	<u>1</u>
<u>Spain</u>	<u>4</u>
<u>Greece</u>	<u>1</u>
<u>Jamaica</u>	<u>4</u>
<u>Poland</u>	<u>6</u>
<u>Honduras</u>	<u>2</u>
<u>Brazil</u>	<u>5</u>
<u>Viet Nam</u>	<u>1</u>
<u>Guam</u>	<u>2</u>
<u>Denmark</u>	<u>1</u>
<u>Switzerland</u>	<u>1</u>
<u>Singapore</u>	<u>1</u>

<u>birthCountry</u>	<u>no of births</u>
<u>Belgium</u>	1
<u>Peru</u>	1
<u>Belize</u>	1
<u>Indonesia</u>	1
<u>Finland</u>	1
<u>Lithuania</u>	1
<u>South Africa</u>	2
<u>Slovakia</u>	1
<u>Saudi Arabia</u>	2
<u>Portugal</u>	1
<u>Latvia</u>	1

## Best Baseball Players

- This question uses `lahmansbaseballdb.batting`, `lahmansbaseballdb.pitching` and `lahmansbaseballdb.people`.
- These query computes performance metrics:
  - Batting:
    - On-base percentage: OBP is  $(\text{sum}(h) + \text{sum}(BB)) / (\text{sum}(ab) + \text{sum}(BB))$ .
    - Slugging percentage: SLG is

$$\frac{(\text{sum}(h) - \text{sum}('2b') - \text{sum}('3b') - \text{sum}(hr)) + 2 * \text{sum}('2b') + 3 * \text{sum}('3b') + 4 * \text{sum}(hr))}{\text{sum}(ab)} \quad (3)$$



- On-base percentage plus slugging: OPS is is  $(obp + slg)$ .
- Pitching:
  - total wins is  $sum(w)$ .
  - total losses is  $sum(l)$ .
  - win percentage is  $sum(w)/(sum(w) + sum(l))$ .
- Professor Ferguson has two criteria for someone being a great baseball player.
  - Batting:
    - Total number of  $ab \geq 3000$ .
    - OPS: Career  $OPS \geq 1.000$
  - Pitching:
    - $(sum(w) + sum(l)) \geq 200$ .
    - $win\_percentage \geq 0.70$  or  $sum(w) \geq 300$ .
- This is one of the rare cases where Prof. Ferguson will provide the answer. So, please produce the table below. Some notes:
  - $great\_because$  is either  $Pitcher$  or  $Batter$  based on whether the player matched the batting or pitching criteria.
  - The values from  $batting$  are  $None$  if the player did not qualify based on batting.
  - The values from  $pitching$  are  $None$  if the player did not qualify on pitching.
- There is a CSV file in the directory with my sample answer. The columns are:
  - $playerid$
  - $nameLast$  from  $People$
  - $nameFirst$  from  $People$ .
  - $great\_because$
  - $slg$
  - $obp$
  - $op\_slg$  is on-base percentage plus slugging.
  - $total\_abs$  is total at bats.
  - $total\_w$  is the total wins.

- total\_l is the total loses.
  - total\_d is some of total\_w and total\_l.
  - win\_percentage is the winning percentage.
- **Note:** Since I saved the query to Pandas and then read it back it,
    - All number become floating point. Your answer uses SQL and the values should be integers except for win\_percentage.
    - Pandas converts NULL to NaN. In your SQL result, you will have None instead of NaN.

```
In [45]: bb_greats = pd.read_csv('./baseball_greats.csv')
bb_greats
```

Out[45]:

	<b>playerid</b>	<b>nameLast</b>	<b>nameFirst</b>	<b>great because</b>	<b>slg</b>	<b>obp</b>	<b>obp slg</b>	<b>total abs</b>	<b>total w</b>	<b>total l</b>	<b>total d</b>	<b>win percentage</b>
<b>0</b>	<a href="#">alexape01</a>	<a href="#">Alexander</a>	<a href="#">Pete</a>	<a href="#">Pitcher</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">373.0</a>	<a href="#">208.0</a>	<a href="#">581.0</a>	<a href="#">0.6420</a>
<b>1</b>	<a href="#">bondsba01</a>	<a href="#">Bonds</a>	<a href="#">Barry</a>	<a href="#">Batter</a>	<a href="#">0.6069</a>	<a href="#">0.4428</a>	<a href="#">1.0497</a>	<a href="#">9847.0</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>
<b>2</b>	<a href="#">carltst01</a>	<a href="#">Carlton</a>	<a href="#">Steve</a>	<a href="#">Pitcher</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">329.0</a>	<a href="#">244.0</a>	<a href="#">573.0</a>	<a href="#">0.5742</a>
<b>3</b>	<a href="#">carutbo01</a>	<a href="#">Caruthers</a>	<a href="#">Bob</a>	<a href="#">Pitcher</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">218.0</a>	<a href="#">99.0</a>	<a href="#">317.0</a>	<a href="#">0.6877</a>
<b>4</b>	<a href="#">clarkjo01</a>	<a href="#">Clarkson</a>	<a href="#">John</a>	<a href="#">Pitcher</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">328.0</a>	<a href="#">178.0</a>	<a href="#">506.0</a>	<a href="#">0.6482</a>
<b>5</b>	<a href="#">clemero02</a>	<a href="#">Clemens</a>	<a href="#">Roger</a>	<a href="#">Pitcher</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">354.0</a>	<a href="#">184.0</a>	<a href="#">538.0</a>	<a href="#">0.6580</a>
<b>6</b>	<a href="#">fordwh01</a>	<a href="#">Ford</a>	<a href="#">Whitey</a>	<a href="#">Pitcher</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">236.0</a>	<a href="#">106.0</a>	<a href="#">342.0</a>	<a href="#">0.6901</a>
<b>7</b>	<a href="#">foutzda01</a>	<a href="#">Foutz</a>	<a href="#">Dave</a>	<a href="#">Pitcher</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">147.0</a>	<a href="#">66.0</a>	<a href="#">213.0</a>	<a href="#">0.6901</a>
<b>8</b>	<a href="#">foxxji01</a>	<a href="#">Foxy</a>	<a href="#">Jimmie</a>	<a href="#">Batter</a>	<a href="#">0.6093</a>	<a href="#">0.4275</a>	<a href="#">1.0368</a>	<a href="#">8134.0</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>
<b>9</b>	<a href="#">galvipu01</a>	<a href="#">Galvin</a>	<a href="#">Pud</a>	<a href="#">Pitcher</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">365.0</a>	<a href="#">310.0</a>	<a href="#">675.0</a>	<a href="#">0.5407</a>
<b>10</b>	<a href="#">gehrilo01</a>	<a href="#">Gehrig</a>	<a href="#">Lou</a>	<a href="#">Batter</a>	<a href="#">0.6324</a>	<a href="#">0.4447</a>	<a href="#">1.0772</a>	<a href="#">8001.0</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>
<b>11</b>	<a href="#">glavito02</a>	<a href="#">Glavine</a>	<a href="#">Tom</a>	<a href="#">Pitcher</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">305.0</a>	<a href="#">203.0</a>	<a href="#">508.0</a>	<a href="#">0.6004</a>
<b>12</b>	<a href="#">greenha01</a>	<a href="#">Greenberg</a>	<a href="#">Hank</a>	<a href="#">Batter</a>	<a href="#">0.6050</a>	<a href="#">0.4103</a>	<a href="#">1.0153</a>	<a href="#">5193.0</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>
<b>13</b>	<a href="#">grovele01</a>	<a href="#">Grove</a>	<a href="#">Lefty</a>	<a href="#">Pitcher</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">300.0</a>	<a href="#">141.0</a>	<a href="#">441.0</a>	<a href="#">0.6803</a>
<b>14</b>	<a href="#">hornsro01</a>	<a href="#">Hornsby</a>	<a href="#">Rogers</a>	<a href="#">Batter</a>	<a href="#">0.5765</a>	<a href="#">0.4308</a>	<a href="#">1.0073</a>	<a href="#">8173.0</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>
<b>15</b>	<a href="#">johnsra05</a>	<a href="#">Johnson</a>	<a href="#">Randy</a>	<a href="#">Pitcher</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">303.0</a>	<a href="#">166.0</a>	<a href="#">469.0</a>	<a href="#">0.6461</a>
<b>16</b>	<a href="#">johnswa01</a>	<a href="#">Johnson</a>	<a href="#">Walter</a>	<a href="#">Pitcher</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">417.0</a>	<a href="#">279.0</a>	<a href="#">696.0</a>	<a href="#">0.5991</a>
<b>17</b>	<a href="#">keefeti01</a>	<a href="#">Keefe</a>	<a href="#">Tim</a>	<a href="#">Pitcher</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">342.0</a>	<a href="#">225.0</a>	<a href="#">567.0</a>	<a href="#">0.6032</a>
<b>18</b>	<a href="#">kershcl01</a>	<a href="#">Kershaw</a>	<a href="#">Clayton</a>	<a href="#">Pitcher</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">185.0</a>	<a href="#">84.0</a>	<a href="#">269.0</a>	<a href="#">0.6877</a>
<b>19</b>	<a href="#">maddugr01</a>	<a href="#">Maddux</a>	<a href="#">Greg</a>	<a href="#">Pitcher</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">355.0</a>	<a href="#">227.0</a>	<a href="#">582.0</a>	<a href="#">0.6100</a>
<b>20</b>	<a href="#">martipe02</a>	<a href="#">Martinez</a>	<a href="#">Pedro</a>	<a href="#">Pitcher</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">219.0</a>	<a href="#">100.0</a>	<a href="#">319.0</a>	<a href="#">0.6865</a>
<b>21</b>	<a href="#">mathech01</a>	<a href="#">Mathewson</a>	<a href="#">Christy</a>	<a href="#">Pitcher</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">373.0</a>	<a href="#">188.0</a>	<a href="#">561.0</a>	<a href="#">0.6649</a>
<b>22</b>	<a href="#">nichoki01</a>	<a href="#">Nichols</a>	<a href="#">Kid</a>	<a href="#">Pitcher</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">362.0</a>	<a href="#">208.0</a>	<a href="#">570.0</a>	<a href="#">0.6351</a>
<b>23</b>	<a href="#">niekrph01</a>	<a href="#">Niekro</a>	<a href="#">Phil</a>	<a href="#">Pitcher</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">NaN</a>	<a href="#">318.0</a>	<a href="#">274.0</a>	<a href="#">592.0</a>	<a href="#">0.5372</a>

	<u>playerid</u>	<u>nameLast</u>	<u>nameFirst</u>	<u>great because</u>	<u>slg</u>	<u>obp</u>	<u>obp slg</u>	<u>total abs</u>	<u>total w</u>	<u>total l</u>	<u>total d</u>	<u>win percentage</u>
24	perryga01	Perry	Gaylord	Pitcher	NaN	NaN	NaN	NaN	314.0	265.0	579.0	0.5423
25	planked01	Plank	Eddie	Pitcher	NaN	NaN	NaN	NaN	326.0	194.0	520.0	0.6269
26	radboch01	Radbourn	Old Hoss	Pitcher	NaN	NaN	NaN	NaN	310.0	194.0	504.0	0.6151
27	ruthba01	Ruth	Babe	Batter	0.6898	0.4718	1.1616	8398.0	NaN	NaN	NaN	NaN
28	ryanno01	Ryan	Nolan	Pitcher	NaN	NaN	NaN	NaN	324.0	292.0	616.0	0.5260
29	seaveto01	Seaver	Tom	Pitcher	NaN	NaN	NaN	NaN	311.0	205.0	516.0	0.6027
30	spahnwa01	Spahn	Warren	Pitcher	NaN	NaN	NaN	NaN	363.0	245.0	608.0	0.5970
31	spaldal01	Spalding	Al	Pitcher	NaN	NaN	NaN	NaN	252.0	65.0	317.0	0.7950
32	suttodo01	Sutton	Don	Pitcher	NaN	NaN	NaN	NaN	324.0	256.0	580.0	0.5586
33	welchmi01	Welch	Mickey	Pitcher	NaN	NaN	NaN	NaN	307.0	210.0	517.0	0.5938
34	willite01	Williams	Ted	Batter	0.6338	0.4806	1.1144	7706.0	NaN	NaN	NaN	NaN
35	wynnea01	Wynn	Early	Pitcher	NaN	NaN	NaN	NaN	300.0	244.0	544.0	0.5515
36	youngcy01	Young	Cy	Pitcher	NaN	NaN	NaN	NaN	511.0	315.0	826.0	0.6186

Answer

In [55]: `%%sql`

```
USE lahmans hwl;
```

```
WITH BattingStats AS (SELECT playerid,
                             (sum(H) + sum(BB)) / (sum(AB) + sum(BB)) AS obp,
                             (sum(H) - sum(2B) - sum(3B) - sum(HR) + 2 * sum(2B) + 3 * sum(3B) + 4 * sum(HR)) /
                             sum(AB) AS slg,
                             ((sum(H) + sum(BB)) / (sum(AB) + sum(BB))) +
                             ((sum(H) - sum(2B) - sum(3B) - sum(HR) + 2 * sum(2B) + 3 * sum(3B) + 4 * sum(HR)) /
                             sum(AB)) AS obp_slg,
                             sum(AB) AS total abs
                             FROM Batting
                             GROUP BY playerid
                             HAVING total abs >= 3000
                             AND obp_slg >= 1.000),
```

```

PitchingStat AS (SELECT playerid,
                    sum(W) as total w,
                    sum(L) as total l,
                    sum(W) / (sum(W) + sum(L)) AS win percentage,
                    sum(W) + sum(L) AS total d
                  FROM Pitching
                  GROUP BY playerid
                  HAVING total d >= 200
                  AND (win percentage >= 0.70 OR total w >= 300))

SELECT People.playerid,
       People.nameLast,
       People.nameFirst,
       CASE
         WHEN BattingStats.playerid IS NOT NULL THEN 'Batter'
         WHEN PitchingStat.playerid IS NOT NULL THEN 'Pitcher'
       END AS great because,
       BattingStats.slg AS slg,
       BattingStats.obp AS obp,
       BattingStats.obp slg AS obp slg,
       BattingStats.total abs AS total abs,
       PitchingStat.total w AS total w,
       PitchingStat.total l AS total l,
       PitchingStat.total d AS total d,
       PitchingStat.win percentage AS win percentage
FROM People
     LEFT JOIN BattingStats
           ON People.playerid = BattingStats.playerid
     LEFT JOIN PitchingStat
           ON People.playerid = PitchingStat.playerid
WHERE BattingStats.playerid IS NOT NULL
     OR PitchingStat.playerid IS NOT NULL;

```

```
* mysql+pymysql://root:***@localhost
```

```
0 rows affected.
```

```
32 rows affected.
```

Out[55]:

<u>playerid</u>	<u>nameLast</u>	<u>nameFirst</u>	<u>great</u>	<u>because</u>	<u>slg</u>	<u>obp</u>	<u>obp slg</u>	<u>total abs</u>	<u>total w</u>	<u>total l</u>	<u>total d</u>	<u>win_percentage</u>
<u>alexape01</u>	<u>Alexander</u>	<u>Pete</u>		<u>Pitcher</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>373</u>	<u>208</u>	<u>581</u>	<u>0.6420</u>
<u>bondsba01</u>	<u>Bonds</u>	<u>Barry</u>		<u>Batter</u>	<u>0.6069</u>	<u>0.4428</u>	<u>1.0497</u>	<u>9847</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>
<u>carltst01</u>	<u>Carlton</u>	<u>Steve</u>		<u>Pitcher</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>329</u>	<u>244</u>	<u>573</u>	<u>0.5742</u>
<u>clarkjo01</u>	<u>Clarkson</u>	<u>John</u>		<u>Pitcher</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>328</u>	<u>178</u>	<u>506</u>	<u>0.6482</u>
<u>clemero02</u>	<u>Clemens</u>	<u>Roger</u>		<u>Pitcher</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>354</u>	<u>184</u>	<u>538</u>	<u>0.6580</u>
<u>foxxji01</u>	<u>Foxx</u>	<u>Jimmie</u>		<u>Batter</u>	<u>0.6093</u>	<u>0.4275</u>	<u>1.0368</u>	<u>8134</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>
<u>galvipu01</u>	<u>Galvin</u>	<u>Pud</u>		<u>Pitcher</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>365</u>	<u>310</u>	<u>675</u>	<u>0.5407</u>
<u>gehrilo01</u>	<u>Gehrig</u>	<u>Lou</u>		<u>Batter</u>	<u>0.6324</u>	<u>0.4447</u>	<u>1.0772</u>	<u>8001</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>
<u>glavito02</u>	<u>Glavine</u>	<u>Tom</u>		<u>Pitcher</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>305</u>	<u>203</u>	<u>508</u>	<u>0.6004</u>
<u>greenha01</u>	<u>Greenberg</u>	<u>Hank</u>		<u>Batter</u>	<u>0.6050</u>	<u>0.4103</u>	<u>1.0153</u>	<u>5193</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>
<u>grovele01</u>	<u>Grove</u>	<u>Lefty</u>		<u>Pitcher</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>300</u>	<u>141</u>	<u>441</u>	<u>0.6803</u>
<u>hornsro01</u>	<u>Hornsby</u>	<u>Rogers</u>		<u>Batter</u>	<u>0.5765</u>	<u>0.4308</u>	<u>1.0073</u>	<u>8173</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>
<u>johnsra05</u>	<u>Johnson</u>	<u>Randy</u>		<u>Pitcher</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>303</u>	<u>166</u>	<u>469</u>	<u>0.6461</u>
<u>johnswa01</u>	<u>Johnson</u>	<u>Walter</u>		<u>Pitcher</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>417</u>	<u>279</u>	<u>696</u>	<u>0.5991</u>
<u>keefeti01</u>	<u>Keefe</u>	<u>Tim</u>		<u>Pitcher</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>342</u>	<u>225</u>	<u>567</u>	<u>0.6032</u>
<u>maddugr01</u>	<u>Maddux</u>	<u>Greg</u>		<u>Pitcher</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>355</u>	<u>227</u>	<u>582</u>	<u>0.6100</u>
<u>mathech01</u>	<u>Mathewson</u>	<u>Christy</u>		<u>Pitcher</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>373</u>	<u>188</u>	<u>561</u>	<u>0.6649</u>
<u>nichoki01</u>	<u>Nichols</u>	<u>Kid</u>		<u>Pitcher</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>362</u>	<u>208</u>	<u>570</u>	<u>0.6351</u>
<u>niekrph01</u>	<u>Niekro</u>	<u>Phil</u>		<u>Pitcher</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>318</u>	<u>274</u>	<u>592</u>	<u>0.5372</u>
<u>perryga01</u>	<u>Perry</u>	<u>Gaylord</u>		<u>Pitcher</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>314</u>	<u>265</u>	<u>579</u>	<u>0.5423</u>
<u>planked01</u>	<u>Plank</u>	<u>Eddie</u>		<u>Pitcher</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>326</u>	<u>194</u>	<u>520</u>	<u>0.6269</u>
<u>radboch01</u>	<u>Radbourn</u>	<u>Old Hoss</u>		<u>Pitcher</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>310</u>	<u>194</u>	<u>504</u>	<u>0.6151</u>
<u>ruthba01</u>	<u>Ruth</u>	<u>Babe</u>		<u>Batter</u>	<u>0.6898</u>	<u>0.4718</u>	<u>1.1616</u>	<u>8398</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>
<u>ryanno01</u>	<u>Ryan</u>	<u>Nolan</u>		<u>Pitcher</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>324</u>	<u>292</u>	<u>616</u>	<u>0.5260</u>

<u>playerid</u>	<u>nameLast</u>	<u>nameFirst</u>	<u>great because</u>	<u>slg</u>	<u>obp</u>	<u>obp slg</u>	<u>total abs</u>	<u>total w</u>	<u>total l</u>	<u>total d</u>	<u>win percentage</u>
<u>seaveto01</u>	<u>Seaver</u>	<u>Tom</u>	<u>Pitcher</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>311</u>	<u>205</u>	<u>516</u>	<u>0.6027</u>
<u>spahnwa01</u>	<u>Spahn</u>	<u>Warren</u>	<u>Pitcher</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>363</u>	<u>245</u>	<u>608</u>	<u>0.5970</u>
<u>spaldal01</u>	<u>Spalding</u>	<u>Al</u>	<u>Pitcher</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>252</u>	<u>65</u>	<u>317</u>	<u>0.7950</u>
<u>suttodo01</u>	<u>Sutton</u>	<u>Don</u>	<u>Pitcher</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>324</u>	<u>256</u>	<u>580</u>	<u>0.5586</u>
<u>welchmi01</u>	<u>Welch</u>	<u>Mickey</u>	<u>Pitcher</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>307</u>	<u>210</u>	<u>517</u>	<u>0.5938</u>
<u>willite01</u>	<u>Williams</u>	<u>Ted</u>	<u>Batter</u>	<u>0.6338</u>	<u>0.4806</u>	<u>1.1144</u>	<u>7706</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>
<u>wynnea01</u>	<u>Wynn</u>	<u>Early</u>	<u>Pitcher</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>300</u>	<u>244</u>	<u>544</u>	<u>0.5515</u>
<u>youngcy01</u>	<u>Young</u>	<u>Cy</u>	<u>Pitcher</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>None</u>	<u>511</u>	<u>315</u>	<u>826</u>	<u>0.6186</u>

In [ ]: \_

## Data and Schema Cleanup

### Part 1: Countries and Cities

- There is a file `worldcities.csv` in the same folder as this notebook.
- In the following code cell, use Pandas to:
  - Read the CSV file into a Data Frame.
  - Convert the Data Frame to contain only the following columns:
    - `city`
    - `city_ascii`
    - `lat`
    - `lng`
    - `country`
    - `iso2`

- `iso3`
- `id`
- Write the data to the table `worldcities` in the schema `F23W4111Midterm`.
- Use the SQL after the code cell to display part of your new table.

**Note:** In lecture examples and HW, you have seen how to use Pandas to read CSV files and create tables.

### Answer

In [7]: `%%sql`

```
DROP SCHEMA IF EXISTS F23W4111Midterm;
CREATE SCHEMA F23W4111Midterm;
USE F23W4111Midterm;
```

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

Out[7]: `[]`

```
In [8]: worldcities = pd.read_csv('./worldcities.csv')
worldcities = worldcities[['city', 'city_ascii', 'lat', 'lng', 'country', 'iso2', 'iso3', 'id']]
worldcities.to_sql(
    'worldcities', schema="F23W4111Midterm",
    con=sql_engine,
    if_exists="replace",
    index=False
)
```

Out[8]: `41001`

- Display data.

In [9]: `%sql select * from F23W4111Midterm.worldcities order by city limit 30;`



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

Out[9]:

<u>city</u>	<u>city_ascii</u>	<u>lat</u>	<u>lng</u>	<u>country</u>	<u>iso2</u>	<u>iso3</u>	<u>id</u>
'Adrā	`Adra	33.6	36.515	Syria	SY	SYR	1760640037
'Ajlūn	`Ajlun	32.3325	35.7517	Jordan	JO	JOR	1400775371
'Ajmān	`Ajman	25.3994	55.4797	United Arab Emirates	AE	ARE	1784337875
'Akko	`Akko	32.9261	35.0839	Israel	IL	ISR	1376781950
'Alavīcheh	`Alavicheh	33.0528	51.0825	Iran	IR	IRN	1364605877
'Amrān	`Amran	15.6594	43.9439	Yemen	YE	YEM	1887433410
'Āmūdā	`Amuda	37.1042	40.93	Syria	SY	SYR	1760247135
'Anadān	`Anadan	36.2936	37.0444	Syria	SY	SYR	1760993442
'Assāl al Ward	`Assal al Ward	33.8658	36.4133	Syria	SY	SYR	1760181042
'Ataq	`Ataq	14.55	46.8	Yemen	YE	YEM	1887172893
'Ayn 'Īsā	`Ayn `Isa	36.3858	38.8472	Syria	SY	SYR	1760078370
'Ibrī	`Ibri	23.2254	56.517	Oman	OM	OMN	1512077267
'Ain Abessa	'Ain Abessa	36.3	5.295	Algeria	DZ	DZA	1012074116
'Ain Arnat	'Ain Arnat	36.1833	5.3167	Algeria	DZ	DZA	1012453452
'Ain Azel	'Ain Azel	35.8433	5.5219	Algeria	DZ	DZA	1012746080
'Ain el Hammam	'Ain el Hammam	36.5647	4.3061	Algeria	DZ	DZA	1012595495
'Ain Leuh	'Ain Leuh	33.2833	-5.3833	Morocco	MA	MAR	1504668626
'Ain Roua	'Ain Roua	36.3344	5.1806	Algeria	DZ	DZA	1012529757
'Ali Ben Sliman	'Ali Ben Sliman	31.9053	-7.2144	Morocco	MA	MAR	1504127885
'Ayn Bni Mathar	'Ayn Bni Mathar	34.0889	-2.0247	Morocco	MA	MAR	1504845272
's-Hertogenbosch	's-Hertogenbosch	51.6833	5.3167	Netherlands	NL	NLD	1528012333
A Coruña	A Coruna	43.3713	-8.4188	Spain	ES	ESP	1724417375
Aachen	Aachen	50.7762	6.0838	Germany	DE	DEU	1276805572
Aadorf	Aadorf	47.4939	8.8975	Switzerland	CH	CHE	1756022542

<u>city</u>	<u>city_ascii</u>	<u>lat</u>	<u>lng</u>	<u>country</u>	<u>iso2</u>	<u>iso3</u>	<u>id</u>
<u>Aalborg</u>	<u>Aalborg</u>	<u>57.0337</u>	<u>9.9166</u>	<u>Denmark</u>	<u>DK</u>	<u>DNK</u>	<u>1208789278</u>
<u>Aalen</u>	<u>Aalen</u>	<u>48.8372</u>	<u>10.0936</u>	<u>Germany</u>	<u>DE</u>	<u>DEU</u>	<u>1276757787</u>
<u>Aalsmeer</u>	<u>Aalsmeer</u>	<u>52.2639</u>	<u>4.7625</u>	<u>Netherlands</u>	<u>NL</u>	<u>NLD</u>	<u>1528899853</u>
<u>Aalst</u>	<u>Aalst</u>	<u>50.9333</u>	<u>4.0333</u>	<u>Belgium</u>	<u>BE</u>	<u>BEL</u>	<u>1056695813</u>
<u>Aalten</u>	<u>Aalten</u>	<u>51.925</u>	<u>6.5808</u>	<u>Netherlands</u>	<u>NL</u>	<u>NLD</u>	<u>1528326020</u>
<u>Äänekoski</u>	<u>Aanekoski</u>	<u>62.6042</u>	<u>25.7264</u>	<u>Finland</u>	<u>FI</u>	<u>FIN</u>	<u>1246710490</u>

## Part 2: Modify World City Data

- Having multiple rows that repeat `country`, `iso2`, and `iso3` is poor design.
- Create two new tables:
  - `countries` that contains `country`, `iso2`, and `iso3`.
  - `cities` that contains only the remaining fields.
  - Pick either `iso2` or `iso3` to define a foreign key between the tables.
- Add primary keys, unique keys, select column data types, etc., to define a better schema for the two tables.
- Show your SQL statements for creating and modifying the tables below.
- **Note:** A small number of the ISO2 and ISO3 codes are incorrect and will prevent you from creating keys. You must correct this data and document your changes.
- Show your DDL below.

Answer

```
In [10]: %%sql

USE F23W4111Midterm;

update worldcities
set iso3 = 'ARM'
where iso2 = 'AM';

update worldcities
set iso2 = 'NA'
where iso3 = 'NAM';

create table countries as select distinct country, iso2, iso3 from worldcities;
create table cities as select city, city ascii, lat, lng, iso3, id from worldcities;

alter table countries
    modify country varchar(50) not null;

alter table countries
    modify iso2 varchar(2) not null;

alter table countries
    modify iso3 varchar(3) not null;

alter table countries
    add constraint countries_pk
    primary key (country);

alter table countries
    add constraint countries_pk2
    unique (iso2, iso3);
```

```
* mysql+pymysql://root:***@localhost
0 rows affected.
67 rows affected.
31 rows affected.
237 rows affected.
41001 rows affected.
237 rows affected.
237 rows affected.
237 rows affected.
0 rows affected.
0 rows affected.
```

Out[10]: []

```
In [11]: %%sql

alter table cities
    modify city varchar(50) not null;

alter table cities
    modify city ascii varchar(50) not null;

alter table cities
    modify lat decimal(9, 4) not null;

alter table cities
    modify lng decimal(9, 4) not null;

alter table cities
    modify iso3 varchar(3) not null;

alter table cities
    add constraint cities_pk
    primary key (id);

create index countries_iso3 index
    on countries (iso3);

alter table cities
    add constraint cities_countries_iso3_fk
    foreign key (iso3) references countries (iso3);
```

```
* mysql+pymysql://root:***@localhost
41001 rows affected.
41001 rows affected.
41001 rows affected.
41001 rows affected.
41001 rows affected.
0 rows affected.
0 rows affected.
41001 rows affected.
[]
```

Out[11]:

## Part 3: An Easy Question

- An interesting question. Is there a better SQL type for latitude and longitude than `DOUBLE` ? If you think there is a better type, what would it be? (You do not need to perform any conversions).

We can combine them as `POINT` data type.

## Part 4: Final Create Table Statements

- Use the DataGrip tool to generate final `CREATE TABLE` statements below. You do not need to execute the statements.

Answer

```
create table if not exists f23w4111midterm.countries
(
    country varchar(50) not null
        primary key,
    iso2    varchar(2) not null,
    iso3    varchar(3) not null,
    constraint countries_pk2
        unique (iso2, iso3)
);

create table if not exists f23w4111midterm.cities
(
```

```

    city      varchar(50) not null,
    city_ascii varchar(50) not null,
    lat       decimal(9, 4) not null,
    lng       decimal(9, 4) not null,
    iso3      varchar(3)  not null,
    id        bigint      not null
        primary key,
    constraint cities_countries_iso3 fk
        foreign key (iso3) references f23w4111midterm.countries (iso3)
);

create index countries_iso3_index
    on f23w4111midterm.countries (iso3);

```

## Part 5: Fixing People Table

### Create a Copy People

- Create a table `F23W4111Midterm.people_modified` that has the same schema and data as `People`.

In [12]: `%%sql`

```

USE F23W4111Midterm;

CREATE TABLE IF NOT EXISTS people_modified AS
    SELECT * FROM lahmans.hw1.People;

```

```

* mysql+pymysql://root:***@localhost
0 rows affected.
20370 rows affected.

```

Out[12]: `[]`

### Fixing ``birthCountry``

- The query below indicates that some `birthCountry` entries in `people` do not map to a known country.

In [13]: `%%sql`

```
select distinct birthCountry, count(*) as count
from people modified
where
    birthCountry not in (select country from countries)
group by birthCountry;
```

\* mysql+pymysql://root:\*\*\*@localhost

10 rows affected.

Out[13]: **birthCountry** **count**

<u>USA</u>	<u>17601</u>
<u>D.R.</u>	<u>826</u>
<u>CAN</u>	<u>258</u>
<u>P.R.</u>	<u>276</u>
<u>Bahamas</u>	<u>7</u>
<u>South Korea</u>	<u>27</u>
<u>Czech Republic</u>	<u>7</u>
<u>V.I.</u>	<u>14</u>
<u>Viet Nam</u>	<u>1</u>
<u>At Sea</u>	<u>1</u>

- My proposed corrections for birthCountry are:

<u>birthCountry</u>	<u>ISO3</u>	<u>ISO2</u>	<u>Correct Country Name</u>
<u>Bahamas</u>	<u>BHS</u>	<u>BS</u>	<u>Bahamas, The</u>
<u>CAN</u>	<u>CAN</u>	<u>CA</u>	<u>Canada</u>
<u>Czech Republic</u>	<u>CZE</u>	<u>CZ</u>	<u>Czechia</u>
<u>South Korea</u>	<u>KOR</u>	<u>KR</u>	<u>Korea, South</u>
<u>USA</u>	<u>USA</u>	<u>US</u>	<u>United States</u>
<u>Viet Nam</u>	<u>VNM</u>	<u>VN</u>	<u>Vietnam</u>



<u>birthCountry</u>	<u>ISO3</u>	<u>ISO2</u>	<u>Correct Country Name</u>
<u>D.R.</u>	<u>DOM</u>	<u>DO</u>	<u>Dominican Republic</u>
<u>P.R.</u>	<u>PRI</u>	<u>PR</u>	<u>Puerto Rico</u>
<u>V.I.</u>	<u>USA</u>	<u>US</u>	<u>United States</u>
<u>At Sea</u>	<u>NULL</u>	<u>NULL</u>	<u>NULL</u>

- Correct people\_modified , making the following changes:
  1. Add a column birthCountryISO3
  2. Correct the entries for birthCountry
  3. Populate the values for birthCountryISO3
  4. Set up a foreign key relationship from people\_modified to countries
  5. Drop the column birthCountry

#### Answer

- Show your SQL statements for altering the table below.
- Run your queries to show correctly modified table.

```
In [14]: %%sql
alter table people_modified
add birthCountryISO3 VARCHAR(3) null;

update people_modified
set birthCountry = 'Bahamas, The',
birthCountryISO3 = 'BHS'
where birthCountry = 'Bahamas';

update people_modified
set birthCountry = 'Canada',
birthCountryISO3 = 'CAN'
where birthCountry = 'CAN';

update people_modified
set birthCountry = 'Czechia',
```

```
birthCountryISO3 = 'CZE'
where birthCountry = 'Czech Republic';

update people modified
set birthCountry = 'Korea, South',
birthCountryISO3 = 'KOR'
where birthCountry = 'South Korea';

update people modified
set birthCountry = 'United States',
birthCountryISO3 = 'USA'
where birthCountry = 'USA';

update people modified
set birthCountry = 'Vietnam',
birthCountryISO3 = 'VNM'
where birthCountry = 'Viet Nam';

update people modified
set birthCountry = 'Dominican Republic',
birthCountryISO3 = 'DOM'
where birthCountry = 'D.R.';

update people modified
set birthCountry = 'Puerto Rico',
birthCountryISO3 = 'PRI'
where birthCountry = 'P.R.';

update people modified
set birthCountry = 'United States',
birthCountryISO3 = 'USA'
where birthCountry = 'V.I.';

update people modified
set birthCountry = NULL,
birthCountryISO3 = NULL
where birthCountry = 'At Sea';

alter table people modified
add constraint people modified countries iso3 fk
foreign key (birthCountryISO3) references countries (iso3);

alter table people modified
drop column birthCountry;
```

```
* mysql+pymysql://root:***@localhost  
0 rows affected.  
7 rows affected.  
258 rows affected.  
7 rows affected.  
27 rows affected.  
17601 rows affected.  
1 rows affected.  
826 rows affected.  
276 rows affected.  
14 rows affected.  
1 rows affected.  
20370 rows affected.  
0 rows affected.
```

Out[14]: []