Administrative Notes March 2, 2023

- March 3: Assignment 4 is due (individual assignment)
 - The question PDF was updated on both Feb 21 and Feb 28 so please make sure you are looking at the right version
 - Tutorials this week will function as open office hours.
 - Friday tutorial moved online for this week only
- March 3 @ 10PM: Assignment 5 released (programming assignments)
 - If you are working in a group, sign up here.
- March 9 (midterm): During class time. Closed book, closed notes, closed neighbor. More details on Piazza (@38).
- Notify Jessica by email ASAP if you have an exam hardship or conflict

CPSC 368 Databases in Data Science

Working with SQL

Learning Goals

- Evaluate the accuracy of given SQL statements for a specified query.
- Determine whether a given query is correctly translated into SQL or not.
- When given a SQL query and a set of relational instances, determine what the SQL query would produce.
- Determine whether two or more SQL queries would produce the same answer for all legal instances of a relation.

General Checklist

- Is the query valid?
 - Are the aggregate operators used properly?
 - O Are we selecting values that exist given the relations we are working with?
 - Are the attributes that are being selected in the grouping list?
- Is the query selecting for the right attributes?
- Try to come up with an English description of the query. If you had to use 30 seconds to state what this query is trying to do, what would you say? See if this description matches with the problem statement.
- Test the query with some sample relational instances
 - Boundary cases
 - o "Normal" case

Orchestra Database

Person(email, name, age)

 This relation stores anyone who has signed up for our mailing list. Tuples in this relation may not be listed in Purchase.

Show(id, year, month, date, showing, attendanceNumber)

Showing describes whether a show was during morning, afternoon, or evening

Song(composer, title)

SongsPerformed(showID, composer, title)

- showID is a foreign key referring to Show
- composer and title are foreign keys referring to attributes of the same name in Song

Purchase(email, showID, price)

- email is a foreign key referring to the email attribute in Person
- showID is a foreign key referring to Show

Musician(<u>id</u>, name, instrument, position, nationality)

PerformedIn(id, showID)

- id refers to the attribute of the same name in Musician
- showID is a foreign key referring to Show

Query: Find the email addresses of everyone who is in the Person table but has not made a purchase before.

Given SQL:

```
SELECT
          email
                                     SELECT email
          Purchase
FROM
                                      FROM
                                            Person
WHERE
          EXISTS(
                                     WHERE NOT EXISTS(
          SELECT email
                                                    SELECT email
          FROM Person
                                                    FROM Purchase
);
                                      );
```

Query: Find the email addresses of everyone who is in the Person table but has not made a purchase before.

Given SQL:

```
email
SELECT
                                      SELECT email
          Purchase
FROM
                                      FROM
                                             Person
WHERE
           EXISTS(
                                      WHERE NOT EXISTS(
           SELECT email
                                                     SFLFCT email
           FROM Person
                                                     FROM Purchase
);
                                      );
```

- We need to retrieve information about people from the Person table, not from the Purchase table
- We are looking for those who have NOT made a purchase, so NOT EXISTS should be used instead of EXISTS.

Query: Find the total number of shows performed during each year and month. If the query is not valid or correct, fix it.

Given SQL:

SELECT DISTINCT year, month, COUNT(date) FROM Show

Query: Find the total number of shows performed during each year and month. If it is not valid or correct, fix it.

Given SQL:

```
SELECT DISTINCT year, month, COUNT(date) FROM Show
```

- The given SQL statement is valid but not correct because it has has an aggregated column (date) while others are not aggregated (year and month)
- In order to count the number of records for each year and month, we need to group the data by those columns, and then count the number of records in each group

```
SELECT year, month, COUNT(date)
FROM Show
GROUP BY year, month
```

9

Query: Find the show IDs of shows where every attendee is older than the average age of all individuals in the Person table.

Given SQL:

```
CREATE VIEW Age(age) AS
SELECT AVG(age)
FROM Purchase
NATURAL JOIN Person
```

```
SELECT showID
FROM Purchase
NATURAL INNER JOIN Person
GROUP BY showID
HAVING AVG(age) >
          (SELECT AVG(age)
          FROM Person)
```

What are these SQL queries trying to do?

Task #3: Fix

Query: Find the show IDs of shows where every attendee is older than the average age of all individuals in the Person table.

```
SELECT showID
FROM Purchase NATURAL INNER JOIN Person
GROUP BY showID
HAVING MIN(age) > (SELECT AVG(age) FROM Person)
```

Definition of Equivalency

 Two queries are equivalent if they produce the same result across all legal instances

- Consider the two queries below meant to find the name of musicians who performed at least 2 songs in each show.
- Are they equivalent?

```
SELECT DISTINCT name
FROM Musician, PerformedIn, SongsPerformed
GROUP BY name
HAVING COUNT(title) > 1
```

```
SELECT name
FROM Musician NATURAL INNER JOIN PerformedIn NATURAL INNER JOIN
SongsPerformed
GROUP BY name, showID
HAVING COUNT(title) > 1
```

Task #4: Take 2!

- Consider the two queries below meant to find the name of musicians who performed at least 2 songs in each show.
- Are they equivalent?

```
SELECT DISTINCT m.name
FROM Musician m, PerformedIn pi, SongsPerformed sp
WHERE m.id = pi.id AND pi.showID = sp.showID
GROUP BY m.name
HAVING COUNT(sp.title) > 1
```

```
SELECT name
FROM Musician NATURAL INNER JOIN PerformedIn NATURAL INNER JOIN
SongsPerformed
GROUP BY name, showID
HAVING COUNT(title) > 1
```

Come up with test cases.

- Only performed in 1 show but played 2 songs
- Only performed in 1 show but played 1 song
- Performed in 2 shows and played at least two in each show
- Performed in 2 shows and at least one of the shows didn't meet the song quota

PerformedIn

ID name1 Carol2 Jia3 Irena4 Jessica

Musician

<u>ID</u>	<u>showID</u>
1	1
2	2
3	3
3	4
4	3
4	5

SongsPerformed

<u>showID</u>	composer	<u>title</u>
1	А	А
1	В	В
2	С	С
3	D	D
3	Е	Е
3	F	F
4	А	Α
4	В	В
5	С	С

```
SELECT DISTINCT m.name
FROM Musician m, PerformedIn pi, SongsPerformed sp
WHERE m.id = pi.id AND pi.showID = sp.showID
GROUP BY m.name
HAVING COUNT(sp.title) > 1
```

SongsPerformed

<u>showID</u>	<u>composer</u>	<u>title</u>
1	А	А
1	В	В
2	С	С
3	D	D
3	Е	Е
3	F	F
4	А	А
4	В	В
5	С	С

Musician

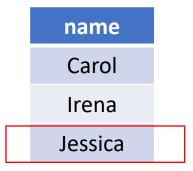
<u>ID</u>	name
1	Carol
2	Jia
3	Irena
4	Jessica

showID
1
2

PerformedIn

Find the name of musicians who performed at least 2 songs in each show.

```
SELECT DISTINCT m.name
FROM Musician m, PerformedIn pi, SongsPerformed sp
WHERE m.id = pi.id AND pi.showID = sp.showID
GROUP BY m.name
HAVING COUNT(sp.title) > 1
```



Find the name of musicians who performed at least 2 songs in each show.

SELECT name
FROM Musician NATURAL INNER JOIN PerformedIn NATURAL INNER JOIN
SongsPerformed
GROUP BY name, showID
SongsPerformed

HAVING COUNT(title) > 1

PerformedIn

ID name1 Carol2 Jia3 Irena4 Jessica

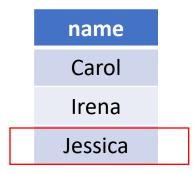
Musician

<u>ID</u>	<u>showID</u>
1	1
2	2
3	3
3	4
4	3
4	5

<u>showID</u>	composer	<u>title</u>
1	А	А
1	В	В
2	С	С
3	D	D
3	Е	Е
3	F	F
4	А	А
4	В	В
5	С	С

Find the name of musicians who performed at least 2 songs in each show.

```
SELECT name
FROM Musician NATURAL INNER JOIN PerformedIn NATURAL INNER JOIN
SongsPerformed
GROUP BY name, showID
HAVING COUNT(title) > 1
```



Task #4: Fix it!

Pick one of the queries and work with the people around you to fix it such that the query can find the name of musicians who performed at least 2 songs in each show.

Task #4: Fix it! Option 1 (amongst others)

Pick one of the queries and work with the people around you to fix it such that the query can find the name of musicians who performed at least 2 songs in each show.

Task #4: Fix it! Option 2 (amongst others)

Pick one of the queries and work with the people around you to fix it such that the query can find the name of musicians who performed at least 2 songs in each show.

```
CREATE VIEW MusiciansWhoPerformedOneSong(id) AS
SELECT m.id
FROM Musician m, PerformedIn pi, SongsPerformed sp
WHERE m.id = pi.id AND pi.showID = sp.showID
GROUP BY pi.showID, m.name
HAVING COUNT(sp.title) = 1
```

```
SELECT m.name
FROM Musician m NATURAL INNER JOIN PerformedIn pi
WHERE m.id NOT IN (SELECT * FROM MusiciansWhoPerformedOneSong)
```

• Consider the query below. If the query is meant to find the title of all songs that have been performed without a pianist, does it accomplish the task?

• Consider the query below. If the query is meant to find the title of all songs that have been performed without a pianist, does it accomplish the task?

- The nested query in the NOT EXISTS clause always returns a row unless there are no pianists in the Musician table.
- This query returns the titles of all songs when there are no pianists at all.
- The problem is that nested query is not CORRELATED with the outer query. In other words, it runs independently of the selected song.

 Consider the query below. If the query is meant to find the title of all songs that have been performed without a pianist, does it accomplish the task?

- When would the inner query be true?
- What relational instances would you use to prove that this query does not work?

• Consider the query below. If the query is meant to find the title of all songs that have been performed without a pianist, does it accomplish the task?

SELECT DISTINCT title
FROM SongsPerformed sp NATURAL INNER JOIN PerformedIn pi
WHERE NOT EXISTS (SELECT *

FROM Musician m
WHERE m.id = pi.id AND

Musician

instrument = "piano")

<u>ID</u>	name	instrument
1	Carol	violin
2	Jia	cello
3	Jessica	piano

SongsPerformed

<u>showID</u>	composer	<u>title</u>
1	А	Α
1	В	В
2	С	С

PerformedIn

<u>ID</u>	<u>showID</u>
1	1
2	1
3	2
1	2

 Consider the query below. If the query is meant to find the title of all songs that have been performed without a pianist, does it accomplish the task?

Result

