Discussion 1

Intro, SQL

Agenda

- I. Welcome!
- II. Getting started
- III. SQL

Welcome!

Course website: http://www.cs186berkeley.net/

Projects are involved coding assignments

- Project 1 is assigned and due Friday 2/2 at 11:59PM

Vitamins are required weekly quizzes that keep you up to date with lectures

- Vitamin 1 due Monday 1/29 at 11:59PM

Midterm 1 will be on 2/21 from 7-9PM Midterm 2 will be on 4/4 from 7-9PM

Accounts

EdStem: https://edstem.org/us/courses/53125/discussion/

Gradescope: See weekly post for add code

GitHub: See Project 0 for GitHub Classroom setup

Email <u>cs186@berkeley.edu</u> if you have any issues

Questions?

SQL

SQL for single tables queries SELECT [DISTINCT] <column list> FROM <table1> [WHERE <predicate>] [GROUP BY <column list>] [HAVING <predicate>] [ORDER BY <column list> [DESC/ASC]] [LIMIT <amount>];

• SQL is declarative - you describe *what* you want in the output, and the DBMS decides *how* it's fetched.

- 1. FROM <table1> which table are we drawing data from
- 2. [WHERE <predicate>] only keep rows where <predicate> is satisfied
- 3. [GROUP BY <column list>] group together rows by value of columns in
 <column list>
- 4. [HAVING <predicate>] only keep groups having <predicate> satisfied
- 5. SELECT <column list> select columns in <column list> to keep
 - a. [DISTINCT] keep only **distinct** rows (filter out duplicates)
- 6. [ORDER BY <column list> [DESC/ASC]] order the output by value of the columns in <column list>, ASCending by default
- 7. [LIMIT <amount>] limit the output to just the first <amount> rows
 - DBMS may execute a query in an equivalent but *different* order
 - For multi-table queries: perform joins with **FROM**

FROM test_table





Take data FROM test_table.

FROM test_table WHERE a > 2







Take data FROM test_table and keep rows WHERE a > 2.

FROM test_table
WHERE a > 2
GROUP BY b



We GROUP BY b, but there may be multiple values of a per group, so we can't use a directly anymore. We can, however, use it with *aggregate functions* (MIN, MAX, SUM, AVERAGE, COUNT). Using aggregate functions without a GROUP BY clause = everything in one group.

FROM test table WHERE a > 2GROUP BY b HAVING COUNT(*) >= 2 b а а

b

Throw away groups that have fewer than 2 rows in them.

Note: COUNT(*) includes NULL values, and COUNT(column) does not include null values.



6

SELECT b AS c FROM test table WHERE a > 2GROUP BY b HAVING COUNT(*) >= 2 b а 3 0 С 5 0 0 4

We use an *alias* here: b AS c selects b, but then calls it c afterwards. We can use this alias in any step *after* this one (so not in SELECT, WHERE, GROUP BY, HAVING).

```
Logical Processing Order
```

SELECT DISTINCT b AS c
FROM test_table
WHERE a > 2
GROUP BY b
HAVING COUNT(*) >= 2



There are no duplicate rows here so DISTINCT isn't necessary, but if there were any, DISTINCT would remove them. Duplicates are removed by exact match *on the entire row*.

```
Logical Processing Order
```

SELECT DISTINCT b AS c FROM test table WHERE a > 2GROUP BY b HAVING COUNT(*) >= 2 ORDER BY c DESC С С 0 1 0

Sort the output by the columns (just c here) (numerically for integers, lexicographically for strings) in either ASCending (low to high) or DESCending (high to low) order. Default is ASC.

Note: order of output is not guaranteed unless you have an ORDER BY clause.

SELECT DISTINCT b AS c FROM test_table WHERE a > 2GROUP BY b HAVING COUNT(*) >= 2 ORDER BY c DESC LIMIT 1 С С 1 1 0

Return just the first row.

A Note On GROUP BY

- Consider the following **Classes** table and
 - Applying the query: **SELECT** <**TBD**> **FROM Classes GROUP BY Dept**;
 - This is what the table looks like once we perform the GROUP BY

| Dept | Course | Capacity |
|------|--------|----------|
| CS | 186 | 715 |
| CS | 161 | 485 |
| EE | 16B | 580 |
| EE | 105 | 40 |
| DATA | 100 | 1375 |

| Dept | Course | Capacity |
|------|--------|----------|
| CS | 186 | 715 |
| CS | 161 | 485 |
| EE | 16B | 580 |
| EE | 105 | 40 |
| DATA | 100 | 1375 |

A Note On GROUP BY

- The role of the SELECT statement is to **squash** each group into a **single row**
 - In general, SELECT clauses are applied to *each row*, but when grouping is involved they are applied to *each group* Are the following valid?

| Dept | Course | Capacity |
|------|--------|----------|
| CS | 186 | 715 |
| CS | 161 | 485 |
| EE | 16B | 580 |
| EE | 105 | 40 |
| DATA | 100 | 1375 |

SELECT Dept, Course FROM Classes GROUP BY Dept;

| Dept | Course |
|------|--------|
| CS | ? |
| EE | ? |
| DATA | ? |

NO

SELECT Dept, SUM(Capacity) FROM Classes GROUP BY Dept;

| Dept | SUM(Cap.) |
|------|-----------|
| CS | 1200 |
| EE | 620 |
| DATA | 1375 |

YES

The Group By Rules

• Columns can only be **selected** if they are a part of the **group by** clause *or* the selected column is an aggregate (i.e. MAX, MIN, AVG, etc.)

Given the table **Students(<u>sid</u>, age, gpa, name, major)**, which of the following are valid?

| SELECT age FROM Students GROUP BY major; | SELECT age FROM Students GROUP BY major, age; | SELECT MAX(age) FROM Students GROUP BY major; | SELECT age, MIN(gpa) FROM Students GROUP BY name; |
|------------------------------------------------|--------------------------------------------------------|-----------------------------------------------------|---------------------------------------------------------|
|------------------------------------------------|--------------------------------------------------------|-----------------------------------------------------|---------------------------------------------------------|

NO

NO

String Comparison

LIKE: following expression follows SQL specified format

- _: Any single character
- %: Zero, one, or more characters
- Looks for a *perfect* string match

Examples:

- LIKE 'z%' starts with z
- LIKE 'z_' exactly 2 letters, 1st is z
- LIKE '_z%' 2nd letter is a z

following expression follows regex
 format

- .: Any single character
- *: Zero, one, or more of the character preceding the symbol
- ^: Match at start of string (If used outside [])
- Looks for *any* pattern in the string that fits

Examples:

- [~] 'z.*' contains z
- ~ `^z.*' starts with z

Note: ~ cannot be used in SQLite (which Project 1 will be using)

Practice: Single-Table Queries

Find the names of the 5 songs that spent the least weeks in the top 40, ordered from least to most. Break ties by song name in alphabetical order.

Tables

```
Songs
(song_id, song_name, album_id,
weeks_in_top_40)
```

Artists
(artist_id, artist_name,
first_yr_active)

Find the names of the 5 songs that spent the least weeks in the top 40, ordered from least to most. Break ties by song name in alphabetical order.

SELECT song_name

FROM Songs

ORDER BY weeks_in_top_40 ASC, song_name ASC

LIMIT 5

Tables

```
Songs
(song_id, song_name, album_id,
weeks_in_top_40)
```

Artists
(artist_id, artist_name,
first_yr_active)

Find the name and first year active of every artist whose name starts with the letter 'B'.

Tables

Songs
(song_id, song_name, album_id,
weeks_in_top_40)

Artists
(artist_id, artist_name,
first_yr_active)

Find the name and first year active of every artist whose name starts with the letter 'B'.

SELECT artist_name,

first_yr_active

FROM artists

WHERE artist_name LIKE 'B%'

Tables

```
Songs
(song_id, song_name, album_id,
weeks_in_top_40)
```

Artists
(artist_id, artist_name,
first_yr_active)

Find the name and first year active of every artist whose name starts with the letter 'B'.

SELECT artist_name,

first_yr_active

FROM artists

WHERE artist name ~ '^B.*'

Tables

```
Songs
(song_id, song_name, album_id,
weeks_in_top_40)
```

Artists
(artist_id, artist_name,
first_yr_active)

Find the total number of albums released per genre.

Tables

Songs
(song_id, song_name, album_id,
weeks_in_top_40)

Artists
(artist_id, artist_name,
first_yr_active)

Find the total number of albums released per genre.

SELECT genre, COUNT(album_id)

FROM Albums

GROUP BY genre;

Tables

Songs
(song_id, song_name, album_id,
weeks_in_top_40)

Artists
(artist_id, artist_name,
first_yr_active)

Find the total number of albums released per genre. Don't include genres with a count less than 10.

Tables

```
Songs
(song_id, song_name, album_id,
weeks_in_top_40)
```

Artists
(artist_id, artist_name,
first_yr_active)

Find the total number of albums released per genre. Don't include genres with a count less than 10.

SELECT genre, COUNT(*)

FROM Albums

GROUP BY genre

HAVING COUNT(*) >= 10;

Tables

```
Songs
(song_id, song_name, album_id,
weeks_in_top_40)
```

Artists
(artist_id, artist_name,
first_yr_active)

Find the genre for which the most albums were released in the year 2000. Assume there are no ties.

Tables

```
Songs
(song_id, song_name, album_id,
weeks_in_top_40)
```

Artists
(artist_id, artist_name,
first_yr_active)

Find the most popular album genre that is released in the year 2000. Assume there are no ties.

SELECT genre

FROM Albums

WHERE yr_released = 2000

GROUP BY genre

ORDER BY COUNT(*) DESC

LIMIT 1;

Tables

```
Songs
(song_id, song_name, album_id,
weeks_in_top_40)
```

Artists
(artist_id, artist_name,
first_yr_active)

SQL Joins

Join Variants

• The different types of joins determine what we do with rows that don't ever match the "join condition"

SELECT * FROM T1 INNER JOIN T2 Join Condition ON T1.a = T2.a;



Ages

Standing

| Name | Age |
|----------|-----|
| Brian | 20 |
| Lakshya | 22 |
| Kimberly | 22 |
| Kimberly | 18 |

| Name | Year |
|----------|----------|
| Brian | Junior |
| Kimberly | Freshman |
| Ben | Senior |
| | |

Ages

| Name | Age |
|----------|-----|
| Brian | 20 |
| Lakshya | 22 |
| Kimberly | 22 |
| Kimberly | 18 |

Standing

| Name | Year |
|----------|----------|
| Brian | Junior |
| Kimberly | Freshman |
| Ben | Senior |
| | |

Result

| Name | Age | Year |
|-------|-----|--------|
| Brian | 20 | Junior |

Ages

| Name | Age |
|----------|-----|
| Brian | 20 |
| Lakshya | 22 |
| Kimberly | 22 |
| Kimberly | 18 |

Standing

| Name | Year |
|----------|----------|
| Brian | Junior |
| Kimberly | Freshman |
| Ben | Senior |
| | |

Result

| Name | Age | Year |
|-------|-----|--------|
| Brian | 20 | Junior |

Ages

| Name | Age |
|----------|-----|
| Brian | 20 |
| Lakshya | 22 |
| Kimberly | 22 |
| Kimberly | 18 |

Standing

| Name | Year |
|----------|----------|
| Brian | Junior |
| Kimberly | Freshman |
| Ben | Senior |
| | |

Result

| Name | Age | Year |
|-------|-----|--------|
| Brian | 20 | Junior |

Ages

| Name | Age |
|----------|-----|
| Brian | 20 |
| Lakshya | 22 |
| Kimberly | 22 |
| Kimberly | 18 |

Standing

| Name | Year |
|----------|----------|
| Brian | Junior |
| Kimberly | Freshman |
| Ben | Senior |
| | |

Result

| Name | Age | Year |
|-------|-----|--------|
| Brian | 20 | Junior |

Ages

| Name | Age |
|----------|-----|
| Brian | 20 |
| Lakshya | 22 |
| Kimberly | 22 |
| Kimberly | 18 |

Standing

| Name | Year |
|----------|----------|
| Brian | Junior |
| Kimberly | Freshman |
| Ben | Senior |
| | |

Result

| Name | Age | Year |
|-------|-----|--------|
| Brian | 20 | Junior |

Ages

| Name | Age |
|----------|-----|
| Brian | 20 |
| Lakshya | 22 |
| Kimberly | 22 |
| Kimberly | 18 |

Standing

| Name | Year |
|----------|----------|
| Brian | Junior |
| Kimberly | Freshman |
| Ben | Senior |

Result

| Name | Age | Year |
|-------|-----|--------|
| Brian | 20 | Junior |

Ages

| Name | Age |
|----------|-----|
| Brian | 20 |
| Lakshya | 22 |
| Kimberly | 22 |
| Kimberly | 18 |

Standing

| Name | Year |
|----------|----------|
| Brian | Junior |
| Kimberly | Freshman |
| Ben | Senior |
| | |

Result

| Name | Age | Year |
|-------|-----|--------|
| Brian | 20 | Junior |

Ages

| Name | Age |
|----------|-----|
| Brian | 20 |
| Lakshya | 22 |
| Kimberly | 22 |
| Kimberly | 18 |

Standing

| Name | Year |
|----------|----------|
| Brian | Junior |
| Kimberly | Freshman |
| Ben | Senior |
| | |

Result

| Name | Age | Year |
|----------|-----|----------|
| Brian | 20 | Junior |
| Kimberly | 22 | Freshman |

Ages

| Name | Age |
|----------|-----|
| Brian | 20 |
| Lakshya | 22 |
| Kimberly | 22 |
| Kimberly | 18 |

Standing

| Name | Year |
|----------|----------|
| Brian | Junior |
| Kimberly | Freshman |
| Ben | Senior |
| | |

Result

| Name | Age | Year |
|----------|-----|----------|
| Brian | 20 | Junior |
| Kimberly | 22 | Freshman |

Ages

| Name | Age |
|----------|-----|
| Brian | 20 |
| Lakshya | 22 |
| Kimberly | 22 |
| Kimberly | 18 |

Standing

| Name | Year |
|----------|----------|
| Brian | Junior |
| Kimberly | Freshman |
| Ben | Senior |
| | |

Result

| Name | Age | Year |
|----------|-----|----------|
| Brian | 20 | Junior |
| Kimberly | 22 | Freshman |

Ages

| Name | Age |
|----------|-----|
| Brian | 20 |
| Lakshya | 22 |
| Kimberly | 22 |
| Kimberly | 18 |

Standing

| Name | Year |
|----------|----------|
| Brian | Junior |
| Kimberly | Freshman |
| Ben | Senior |
| | |

Result

| Name | Age | Year |
|----------|-----|----------|
| Brian | 20 | Junior |
| Kimberly | 22 | Freshman |
| Kimberly | 18 | Freshman |

Ages

| Name | Age |
|----------|-----|
| Brian | 20 |
| Lakshya | 22 |
| Kimberly | 22 |
| Kimberly | 18 |

Standing

| Year |
|----------|
| Junior |
| Freshman |
| Senior |
| |

Result

| Name | Age | Year |
|----------|-----|----------|
| Brian | 20 | Junior |
| Kimberly | 22 | Freshman |
| Kimberly | 18 | Freshman |

Ages

| Name | Age |
|----------|-----|
| Brian | 20 |
| Lakshya | 22 |
| Kimberly | 22 |
| Kimberly | 18 |

Standing

| Name | Year |
|----------|----------|
| Brian | Junior |
| Kimberly | Freshman |
| Ben | Senior |
| | |

Result

| Name | Age | Year |
|----------|-----|----------|
| Brian | 20 | Junior |
| Kimberly | 22 | Freshman |
| Kimberly | 18 | Freshman |

Left Join, Example

Ages

Standing

| Name | Age |
|----------|-----|
| Brian | 20 |
| Lakshya | 22 |
| Kimberly | 22 |
| Kimberly | 18 |

| Name | Year |
|----------|----------|
| Brian | Junior |
| Kimberly | Freshman |
| Ben | Senior |

Left Join, Example

Ages

| Name | Age |
|----------|-----|
| Brian | 20 |
| Lakshya | 22 |
| Kimberly | 22 |
| Kimberly | 18 |

Standing

| inior |
|---------|
| |
| reshman |
| enior |
| , |

Result

| Name | Age | Year |
|----------|-----|----------|
| Brian | 20 | Junior |
| Kimberly | 22 | Freshman |
| Kimberly | 18 | Freshman |
| Lakshya | 22 | null |

Right Join, Example

Ages

Standing

| Name | Age |
|----------|-----|
| Brian | 20 |
| Lakshya | 22 |
| Kimberly | 22 |
| Kimberly | 18 |

| Name | Year |
|----------|----------|
| Brian | Junior |
| Kimberly | Freshman |
| Ben | Senior |

Right Join, Example

Ages

| Name | Age |
|----------|-----|
| Brian | 20 |
| Lakshya | 22 |
| Kimberly | 22 |
| Kimberly | 18 |

Standing

| Name | Year |
|----------|----------|
| Brian | Junior |
| Kimberly | Freshman |
| Ben | Senior |
| | |

Result

| Name | Age | Year |
|----------|------|----------|
| Brian | 20 | Junior |
| Kimberly | 22 | Freshman |
| Kimberly | 18 | Freshman |
| null | null | Senior |

Full Outer Join, Example

Ages

Standing

| Name | Age |
|----------|-----|
| Brian | 20 |
| Lakshya | 22 |
| Kimberly | 22 |
| Kimberly | 18 |

| Name | Year |
|----------|----------|
| Brian | Junior |
| Kimberly | Freshman |
| Ben | Senior |

Full Outer Join, Example

Ages

| Name | Age |
|----------|-----|
| Brian | 20 |
| Lakshya | 22 |
| Kimberly | 22 |
| Kimberly | 18 |

Standing

| Name | Year | |
|----------|----------|--|
| Brian | Junior | |
| Kimberly | Freshman | |
| Ben | Senior | |
| | Genio | |

Result

| Name | Age | Year |
|----------|------|----------|
| Brian | 20 | Junior |
| Kimberly | 22 | Freshman |
| Kimberly | 18 | Freshman |
| null | null | Senior |
| Lakshya | 22 | null |

Practice: Multi-Table Joins

Find the names of all artists who released a 'country' genre album in 2020.

Tables

Songs
(song_id, song_name, album_id,
weeks_in_top_40)

Artists
(artist_id, artist_name,
first_yr_active)

Find the names of all artists who released a 'country' genre album in 2020.

SELECT artist_name

FROM Artists AS A

INNER JOIN Albums AS B

ON A.artist_id = B.artist_id

WHERE genre = 'country' AND

yr_released = 2020

```
GROUP BY A.artist_id,
```

artist_name;

Tables

```
Songs
(song_id, song_name, album_id,
weeks_in_top_40)
```

Artists
(artist_id, artist_name,
first_yr_active)

Find the name of the album with the song that spend the most weeks in the top 40. Assume there is only one such song.

Tables

```
Songs
(song_id, song_name, album_id,
weeks_in_top_40)
```

Artists
(artist_id, artist_name,
first_yr_active)

Find the name of the album with the song that spend the most weeks in the top 40. Assume there is only one such song.

SELECT album_name

FROM Albums AS A INNER JOIN

Songs AS S

ON A.album_id = S.album_id

ORDER BY weeks_in_top_40 DESC
LIMIT 1;

Tables

```
Songs
(song_id, song_name, album_id,
weeks_in_top_40)
```

Artists
(artist_id, artist_name,
first_yr_active)

Find the artist name and the most weeks one of their songs spent in the top 40 for each artist. Include artists that have not released an album.

Tables

```
Songs
(song_id, song_name, album_id,
weeks_in_top_40)
```

Artists
(artist_id, artist_name,
first_yr_active)

Find the artist name and the most weeks one of their songs spent in the top 40 for each artist. Include artists that have not released an album.

SELECT artist_name, MAX(weeks_in_top_40)

FROM Artists LEFT JOIN

(Songs INNER JOIN Albums ON
Songs.album id = Albums.album id)

ON Artists.artist_id = Albums.artist_id

GROUP BY Artists.artist_id, artist_name;

Tables

```
Songs
(song_id, song_name, album_id,
weeks_in_top_40)
```

Artists
(artist_id, artist_name,
first_yr_active)

Appendix (Tips and Tricks)

- Using an aggregate in **WHERE** is **not** allowed
 - WHERE COUNT(*) > 500 is an invalid query!
- Don't use **HAVING** without **GROUP BY**
 - Just use WHERE instead!
- If **GROUP BY** is used, you can only **SELECT** columns that are **aggregates** OR are columns used to **group by**
- **DISTINCT** removes all **duplicate rows**
 - Watch out for order of operators + logical processing order!
 - If you want to see how many distinct values there are of a column,
 DISTINCT COUNT(X) will not work; use COUNT(DISTINCT X)

Attendance Link

https://cs186berkeley.net/attendance

