


Computational Structures in Data Science

UC Berkeley EECS
Adj. Ass. Prof.
Dr. Gerald Friedland

Lecture #12: SQL

April 22nd, 2019 <http://inst.eecs.berkeley.edu/~cs88>


Btw. CNN criticizing my work...



<https://www.cnn.com/2019/04/19/tech/ai-facial-recognition/index.html>

04/22/19 UCB CS88 Sp19 L12 2

Question

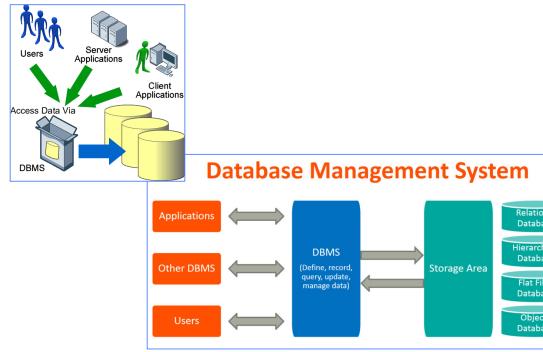
Next lecture...

A) Q&A Session
B) Fun Lecture
C) Review Lecture
D) All of the above



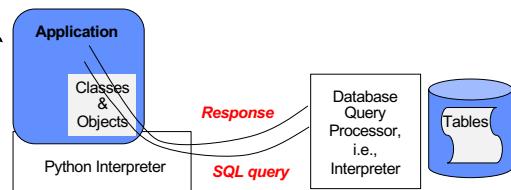
04/15/19 UCB CS88 Sp19 L11


Database Management Systems



04/22/19 UCB CS88 Sp19 L12 4

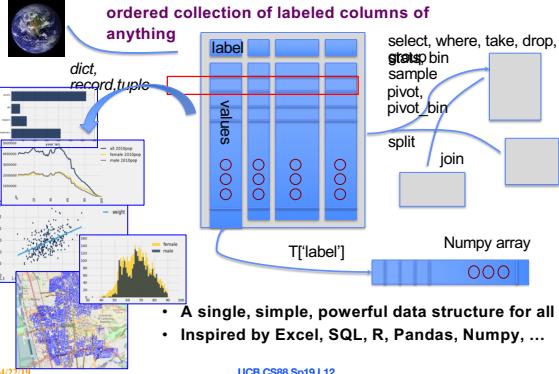
App in program language issues queries to a database interpreter



- The SQL language is represented in query strings delivered to a DB backend.
- Use the techniques learned here to build clean abstractions.
- You have already learned the relational operators!

04/22/19 UCB CS88 Sp19 L12 5


Data 8 Tables



- A single, simple, powerful data structure for all
- Inspired by Excel, SQL, R, Pandas, Numpy, ...

04/22/19 UCB CS88 Sp19 L12 6

Database Management Systems

- DBMS are persistent tables with powerful relational operators
 - Important, heavily used, interesting!
- A table is a collection of records, which are rows that have a value for each column
 - row has a value for each column
 - table has columns and rows
 - column has a name and a type
- Structure Query Language (SQL) is a declarative programming language describing operations on tables

04/22/19 UCB CS88 Sp19 L12 7

SQL

- A declarative language
 - Described what to compute
 - Imperative languages, like python, describe how to compute it
 - Query processor (interpreter) chooses which of many equivalent query plans to execute to perform the SQL statements
- ANSI and ISO standard, but many variants
- select statement creates a new table, either from scratch or by projecting a table
- create table statement gives a global name to a table
- Lots of other statements
 - analyze, delete, explain, insert, replace, update, ...
- The action is in select

04/22/19 UCB CS88 Sp19 L12 8

SQL example

- SQL statements create tables
 - Give it a try with sqlite3 or <http://kripken.github.io/sql.js/GUI/>
 - Each statement ends with ;'

```
culler$ sqlite3
SQLite version 3.9.2 2015-11-02 18:31:45
Enter ".help" for usage hints.
Connected to a transient in-memory database.
Use ".open FILENAME" to reopen on a persistent database.
sqlite> select 38 as latitude, 122 as longitude, 'Berkeley' as name;
38|122|Berkeley
sqlite>
```

04/22/19 UCB CS88 Sp19 L12 9

A Running example from Data 8 Lec 10

```
# An example of creating a Table from a list of rows.
Table([{"Flavor": "strawberry", "Color": "pink", "Price": 3.55},
       {"chocolate", "light brown", 4.75},
       {"chocolate", "dark brown", 5.25},
       {"strawberry", "pink", 5.25},
       {"bubblegum", "pink", 4.75}])
```

Flavor	Color	Price
strawberry	pink	3.55
chocolate	light brown	4.75
chocolate	dark brown	5.25
strawberry	pink	5.25
bubblegum	pink	4.75



culler@CullerMac ~/Classes/CS88-Fel18/ideas/sql> sqlite3 icecream.db
SQLite version 3.13.0 2016-05-18 10:57:30
Enter ".help" for usage hints.
sqlite>

04/22/19 UCB CS88 Sp19 L12 10

select

- Comma-separated list of column descriptions
- Column description is an expression, optionally followed by as and a column name

```
select [expression] as [name], [expression] as [name];...
```

- Selecting literals creates a one-row table

```
select "strawberry" as Flavor, "pink" as Color, 3.55 as Price;
```

- union of select statements is a table containing the union of the rows

```
select "strawberry" as Flavor, "pink" as Color, 3.55 as Price union
select "chocolate", "light brown", 4.75 union
select "chocolate", "dark brown", 5.25 union
select "strawberry", "pink", 5.25 union
select "bubblegum", "pink", 4.75;
```

04/22/19 UCB CS88 Sp19 L12 11

create table

- SQL often used interactively
 - Result of select displayed to the user, but not stored
- Create table statement gives the result a name
 - Like a variable, but for a permanent object

```
create table [name] as [select statement];
```

04/22/19 UCB CS88 Sp19 L12 12

SQL: creating a named table

```
create table cones as
select 1 as ID, "strawberry" as Flavor, "pink" as Color,
3.55 as Price union
select 2, "chocolate", "light brown", 4.75 union
select 3, "chocolate", "dark brown", 5.25 union
select 4, "strawberry", "pink", 5.25 union
select 5, "bubblegum", "pink", 4.75 union
select 6, "chocolate", "dark brown", 5.25;
```

Notice how column names are introduced and implicit later on.

04/22/19

UCB CS88 Sp19 L12

13

Select ...

```
culler@CullerMac ~/Classes/CS88-Fa18/ideas/sql> sqlite3 icecream.db
SQLite version 3.13.0 2016-05-18 10:57:38
Enter ".help" for usage hints.
sqlite> create table cones as
...>   select 1 as ID, "strawberry" as Flavor, "pink" as Color, 3.55 as Price
ce union
...>   select 2, "chocolate", "light brown", 4.75 union
...>   select 3, "chocolate", "dark brown", 5.25 union
...>   select 4, "strawberry", "pink", 5.25 union
...>   select 5, "bubblegum", "pink", 4.75 union
...>   select 6, "chocolate", "dark brown", 5.25;
sqlite> select * from cones;
1|strawberry|pink|3.55
2|chocolate|light brown|4.75
3|chocolate|dark brown|5.25
4|strawberry|pink|5.25
5|bubblegum|pink|4.75
6|chocolate|dark brown|5.25
sqlite> 
```

ID	Flavor	Color	Price
1	strawberry	pink	3.55
2	chocolate	light brown	4.75
3	chocolate	dark brown	5.25
4	strawberry	pink	5.25
5	bubblegum	pink	4.75
6	chocolate	dark brown	5.25

04/22/19

UCB CS88 Sp19 L12

14

Projecting existing tables

- Input table specified by **from** clause
- Subset of rows selected using a **where** clause
- Ordering of the selected rows declared using an **order by** clause

```
select [columns] from [table] where [condition] order by [order];
```

```
select * from cones order by Price;
```

ID	Flavor	Color	Price
1	strawberry	pink	3.55
2	chocolate	light brown	4.75
5	bubblegum	pink	4.75
3	chocolate	dark brown	5.25
4	strawberry	pink	5.25
6	chocolate	dark brown	5.25

04/22/19

UCB CS88 Sp19 L12

15

Projection

```
In [5]: cones.select(['Flavor', 'Price'])
Out[5]:
```

Flavor	Price
strawberry	3.55
chocolate	4.75
chocolate	5.25
strawberry	5.25
bubblegum	4.75
chocolate	5.25

```
sqlite> select Flavor, Price from cones;
Flavor|Price
strawberry|3.55
chocolate|4.75
chocolate|5.25
strawberry|5.25
bubblegum|4.75
chocolate|5.25
```

- Select versus indexing a column?

04/22/19

UCB CS88 Sp19 L12

16

Permanent Data Storage



```
sqlite>.quit
culler@CullerMac ~/Classes/CS88-Fa18/ideas/sql> sqlite3 icecream.db
SQLite version 3.13.0 2016-05-18 10:57:30
Enter ".help" for usage hints.
sqlite>.tables
cones
sqlite> select * from cones where Color is "dark brown";
3|chocolate|dark brown|5.25
6|chocolate|dark brown|5.25
sqlite> 
```

04/22/19

UCB CS88 Sp19 L12

17

Filtering rows - where

- Set of Table records (rows) that satisfy a condition

```
select [columns] from [table] where [condition] order by [order];
```

```
In [5]: cones.select(['Flavor', 'Price'])
Out[5]:
```

Flavor	Price
strawberry	3.55
chocolate	4.75
chocolate	5.25
strawberry	5.25
bubblegum	4.75
chocolate	5.25

```
sqlite> select * from cones where Flavor = "chocolate";
ID|Flavor|Color|Price
2|chocolate|light brown|4.75
3|chocolate|dark brown|5.25
6|chocolate|dark brown|5.25
```

ID	Flavor	Color	Price
3	chocolate	dark brown	5.25
4	strawberry	pink	5.25
6	chocolate	dark brown	5.25

```
SQL:
sqlite> select * from cones where Price > 5;
ID|Flavor|Color|Price
3|chocolate|dark brown|5.25
4|strawberry|pink|5.25
6|chocolate|dark brown|5.25
```

04/22/19

UCB CS88 Sp19 L12

18

SQL Operators for predicate

- use the **WHERE** clause in the SQL statements such as **SELECT**, **UPDATE** and **DELETE** to filter rows that do not meet a specified condition

SQLite understands the following binary operators, in order from highest to lowest precedence:

```
||   *
   /   %
  +   -
<<  >>  &   |
<   <=   >   >=
=   ==   !=   <>  IS   IS NOT  IN   LIKE   GLOB   MATCH   REGEXP
AND
OR
```

Supported unary prefix operators are these:

```
-   +   -   NOT
```

04/22/19

UCB CS88 Sp19 L12



19

Approximate Matching ...

Regular expression matches are so common that they are built in in SQL.

```
sqlite> select * from cones where Flavor like "%berry%";  
Flavor|Color|Price  
strawberry|pink|3.55  
strawberry|pink|5.25  
sqlite>
```

On the other hand, you have the full power of Python to express what you mean.

```
cones.where(cones.apply(lambda x:'berry' in x, 'Flavor'))
```

ID	Flavor	Color	Price
1	strawberry	pink	3.55
4	strawberry	pink	5.25

04/22/19

UCB CS88 Sp19 L12

20

Group and Aggregate

- The **GROUP BY** clause is used to group rows returned by **SELECT statement** into a set of summary rows or groups based on values of columns or expressions.
- Apply an **aggregate function**, such as **SUM**, **AVG**, **MIN**, **MAX** or **COUNT**, to each group to output the summary information.

```
cones.groupby('Flavor')  
Flavor count  
bubblegum 1  
chocolate 3  
strawberry 2
```

```
sqlite> select count(Price), Flavor from cones group by Flavor;  
count|Price|Flavor  
1|bubblegum  
2|chocolate  
2|strawberry
```

```
cones.select(['Flavor', 'Price']).groupby('Flavor').np.mean()  
Flavor Price mean  
bubblegum 4.75  
chocolate 5.0833  
strawberry 4.4
```

```
sqlite> select avg(Price), Flavor from cones group by Flavor;  
avg(Price)|Flavor  
4.75|bubblegum  
5.0|chocolate  
4.4|strawberry
```

04/22/19

UCB CS88 Sp19 L12

21



Unique / Distinct values

```
select DISTINCT [columns] from [table] where [condition] order by [order];
```

```
sqlite> select distinct Flavor, Color from cones;  
strawberry|pink  
chocolate|light brown  
chocolate|dark brown  
bubblegum|pink  
sqlite> 
```

```
In [8]: cones.groupby(['Flavor', 'Color']).drop('count')  
Out[8]:  
Flavor Color  
-----  
bubblegum|pink  
chocolate|dark brown  
chocolate|light brown  
strawberry|pink
```

```
In [17]: np.unique(cones['Flavor'])  
Out[17]: array(['bubblegum', 'chocolate', 'strawberry'], dtype='|S12')
```

- Built in to the language or a composable tool?

04/22/19

UCB CS88 Sp19 L12

22

Joining tables

- Two tables are joined by a comma to yield all combinations of a row from each

```
— select * from sales, cones;
```

```
create table sales as  
select 'Baskin' as Cashier, 1 as TID union  
select 'Baskin', 3 union  
select 'Baskin', 4 union  
select 'Robin', 2 union  
select 'Robin', 5 union  
select 'Robin', 6;
```

```
sales.join('TID', cones, 'ID')  
TID Cashier Flavor Color Price  
1 Baskin strawberry pink 3.55  
2 Robin chocolate light brown 4.75  
3 Baskin chocolate dark brown 5.25  
4 Baskin strawberry pink 5.25  
5 Robin bubblegum pink 4.75  
6 Robin chocolate dark brown 5.25
```

04/22/19

UCB CS88 Sp19 L12



23

Inner Join

```
select * from sales, cones where TID=ID;
```

```
sales.join('TID', cones, 'ID')  
TID Cashier Flavor Color Price  
1 Baskin strawberry pink 3.55  
2 Robin chocolate light brown 4.75  
3 Baskin chocolate dark brown 5.25  
4 Baskin strawberry pink 5.25  
5 Robin bubblegum pink 4.75  
6 Robin chocolate dark brown 5.25
```

04/22/19

UCB CS88 Sp19 L12

24



SQL: using named tables - from

```
select "delicious" as Taste, Flavor, Color from cones
  where Flavor is "chocolate" union
select "other", Flavor, Color from cones
  where Flavor is not "chocolate";
```

```
sqlite> select "delicious" as Taste, Flavor, Color from cones where Flavor is "chocolate" union
...>   ... select "other", Flavor, Color from cones where Flavor is not "chocolate";
Taste|Flavor|Color
delicious|chocolate|dark brown
delicious|chocolate|light brown
other|bubblegum|pink
other|strawberry|pink
sqlite> ||
```

04/22/19

UCB CS88 Sp19 L12

25

Queries within queries

- Any place that a table is named within a select statement, a table could be computed
 - As a sub-query

```
select TID from sales where Cashier is "Baskin";
select * from cones
  where ID in (select TID from sales where Cashier is "Baskin");
sqlite> select * from cones
...>   where ID in (select TID from sales where Cashier is "Baskin");
ID|Flavor|Color|Price
1|strawberry|pink|3.55
3|chocolate|dark brown|5.25
4|strawberry|pink|5.25
```

04/22/19

UCB CS88 Sp19 L12

26

Inserting new records (rows)

```
INSERT INTO table(column1, column2,...)
VALUES (value1, value2,...);
```

```
sqlite> insert into cones(ID, Flavor, Color, Price) values (7, "Vanilla", "White", 3.95);
sqlite> select * from cones;
1|Flavor|Color|Price
1|strawberry|pink|3.55
2|chocolate|light brown|4.75
3|chocolate|dark brown|5.25
4|strawberry|pink|5.25
5|bubblegum|pink|4.75
6|chocolate|dark brown|5.25
7|Vanilla|White|3.95
sqlite> ||
```

ID	Flavor	Color	Price
1	strawberry	pink	3.55
2	chocolate	light brown	4.75
3	chocolate	dark brown	5.25
4	strawberry	pink	5.25
5	bubblegum	pink	4.75
6	chocolate	dark brown	5.25
7	Vanilla	White	3.95

- A database table is typically a shared, durable repository shared by multiple applications

04/22/19

UCB CS88 Sp19 L12

27

Multiple clients of the database

```
Last login: Mon Nov 19 19:42:47 on ttys001
dash@dash-MacBook-Pro:~$ cd /Users/dash/anaconda/envs/datscience/bin
dash@dash-MacBook-Pro:~/anaconda/envs/datscience/bin$ PATH
dash@dash-MacBook-Pro:~/anaconda/envs/datscience/bin$ cd Classes/CS88-Fall2018/ideas/vanilla
dash@dash-MacBook-Pro:~/anaconda/envs/datscience/bin$ sqlite3 icecream.db
SQLite version 3.13.8 2016-08-10 18:57:38
copyright (c) 2016 The SQLite Authors
sqlite> insert into cones(ID, Flavor, Color, Price) values (9, "Fudge", "dark", 7.95);
sqlite> ||
```

ID	Flavor	Color	Price
1	strawberry	pink	3.55
2	chocolate	light brown	4.75
3	chocolate	dark brown	5.25
4	strawberry	pink	5.25
5	bubblegum	pink	4.75
6	chocolate	dark brown	5.25
7	Vanilla	White	3.95
8	Fudge	dark	7.95

- All of the inserts update the common repository

04/22/19

UCB CS88 Sp19 L12

28

SQLite python API

```
In [64]: import sqlite3
In [65]: icecream = sqlite3.connect('icecream.db')
In [66]: icecream.execute('SELECT * FROM cones;')

Out[66]: <sqlite3.Cursor at 0x111127960>

In [67]: icecream.execute('SELECT DISTINCT Flavor FROM cones;').fetchall()
Out[67]: [('strawberry'), ('chocolate'), ('bubblegum')]

In [68]: icecream.execute('SELECT * FROM cones WHERE Flavor is "chocolate";').fetchall()
Out[68]: [(2, 'chocolate', 'light brown', 4.75),
          (3, 'chocolate', 'dark brown', 5.25),
          (6, 'chocolate', 'dark brown', 5.25)]
```

04/22/19

UCB CS88 Sp19 L12

29

Creating DB Abstractions

```
class SQL_Table(Table):
    """ Extend Table class with methods to read/write a Table
        from/to a table in a SQLite3 database.
    """
    @classmethod
    def read(cls, filepath, table, verbose=False):
        """Create a SQL_Table by reading a table from a SQL database."""
        dbconn = sqlite3.connect(filepath,
                               detect_types=sqlite3.PARSE_COLNAMES)

        col_names = sqlcol_names(dbconn, table)
        rows = sqlexec(dbconn, 'SELECT * from %s;' % table, verbose).fetchall()
        dbconn.close()
        return cls(col_names).with_rows(rows)
```

04/22/19

UCB CS88 Sp19 L12

30

DB Abstraction (cont)

```
class SQL_Table(Table):
    ...
    def write(self, filepath, table, verbose=False, overwrite=True):
        """Write a Table into a SQL database as a SQL table."""

        dbconn = sqlite3.connect(filepath)
        # Create table and insert each row
        cols = build_list(self.labels)
        sqlexec(dbconn, "CREATE TABLE %s %s;" % (table, cols), verbose)
        for row in self.rows:
            sqlexec(dbconn, 'INSERT INTO %s VALUES %s;' % (table, tuple(row)))
        dbconn.commit()
        dbconn.close()

    @classmethod
    def cast(cls, table):
        """Return a SQL Table version of a Table."""
        return cls().with_columns(zip(table.labels, table.columns))
```

04/22/19

UCB CS88 Sp19 L12

31

Summary – Part 1

```
SELECT <col spec> FROM <table spec> WHERE <cond spec>
    GROUP BY <group spec> ORDER BY <order spec>;
```



```
INSERT INTO table(column1, column2,...)
    VALUES (value1, value2,...);
```



```
CREATE TABLE name (<columns>);
```



```
CREATE TABLE name AS <select statement>;
```



```
DROP TABLE name;
```

04/22/19

UCB CS88 Sp19 L12

32

Summary

- SQL a declarative programming language on relational tables
 - largely familiar to you from data8
 - create, select, where, order, group by, join
- Databases are accessed through Applications
 - e.g., all modern web apps have Database backend
 - Queries are issued through API
 - » Be careful about app corrupting the database
- Data analytics tend to draw database into memory and operate on it as a data structure
 - e.g., Tables
- More in lab

04/22/19

UCB CS88 Sp19 L12

33

Solutions for the Wandering Mind

- 1) Adding two n-bit integers, how many bits can the result have?
Maximally $n+1$ bits (1 overflow bit).
 - 2) Multiplying two n bit integers, how many bits can the result have?
 $\log a \cdot b = \log a + \log b \Rightarrow \log n \cdot n = 2 \cdot \log n$. $2 \cdot n$ bits.
- Assume:
- a) Exceptions don't exist
 - b) We only reserve 8bit for an integer variable (0-255)
- Questions:
- 1) What would be the result of an addition 255+255?
(255+255) modulo 255=0. One entire summand is lost!
 - 2) What would be the result of a multiplication 255*255?
(255*255) modulo 255=0. The error is: 64770!

04/22/19

UCB CS88 Sp19 L12

34

Solutions for the Wandering Mind

3) Assume l additions of 8bit integers into the same 8bit variable. Can you formulate the maximum error that can occur as a function of l?

Each addition can maximally lose the entire 8bits. The error accumulates. This is, each addition loses maximally 8bits. So it's $l \cdot 8$ bits of loss. So the maximum error 255^l .

For multiplication it's 255^l . This is, exponential error!
(Also compare Lyapunov Exponent in physics)

This is taught in CS61C

04/22/19

UCB CS88 Sp19 L12

35