

STATS 507

Data Analysis in Python

Week 9: Databases (SQL) & PyTorch (part 1)

Adapted from slides by Keith Levin and C. Budak

Recap: HTML, XML and JSON

Each provided a different (though similar) way of storing data

Key motivation of JSON (and, sort of, HTML and XML): self-description

But we saw that JSON could get quite unwieldy quite quickly...

Example of a more complicated JSON object

```
1 complex_json_string="{
2   "id": "0001",
3   "type": "donut",
4   "name": "Cake",
5   "ppu": 0.55,
6   "batters":
7     {
8       "batter":
9         [
10          { "id": "1001", "type": "Regular" },
11          { "id": "1002", "type": "Chocolate" },
12          { "id": "1003", "type": "Blueberry" },
13          { "id": "1004", "type": "Devil's Food" }
14        ]
15      },
16   "topping":
17     [
18       { "id": "5001", "type": "None" },
19       { "id": "5002", "type": "Glazed" },
20       { "id": "5005", "type": "Sugar" },
21       { "id": "5007", "type": "Powdered Sugar" },
22       { "id": "5006", "type": "Chocolate with Sprinkles" },
23       { "id": "5003", "type": "Chocolate" },
24       { "id": "5004", "type": "Maple" }
25     ]
26 }"
```

What if I have hundreds of different kinds of cakes or donuts? The nestedness of JSON objects makes them a little complicated. Generally, JSON is good for delivering (small amounts of) data, but for storing and manipulating large, complicated collections of data, there are better tools, namely databases.

Note: there are also security and software engineering reasons to prefer databases over JSON for storing data, but that's beyond the scope of our course.

Why use a database?

Database software hides the problem of actually handling data

As we'll see in a few slides, this is a complicated thing to do!

Indexing, journaling, archiving handled automatically

Allow fast, concurrent (i.e., multiple users) access to data

ACID transactions (more on this in a few slides)

Access over the web

DBs can be run, e.g., on a server

Again, JSON/XML/HTML/etc good for delivering data, DBs good for storing

Databases (DBs)

Information, organized so as to make retrieval fast and efficient

Examples: Census information, product inventory, library catalogue

This course: **relational databases**

https://en.wikipedia.org/wiki/Relational_database

So-named because they capture relations between entities

In existence since the 1970s, and still the dominant model in use today

Outside the scope of this course: other models (e.g., object-oriented)

https://en.wikipedia.org/wiki/Database_model

Textbook: *Database System Concepts* by Silberschatz, Korth and Sudarshan.

Relational DBs: pros and cons

Pros:

- Natural for the vast majority of applications
- Numerous tools for managing and querying

Cons:

- Not well-suited to some data (e.g., networks, unstructured text)
- Fixed schema (i.e., hard to add columns)
- Expensive to maintain when data gets large (e.g., many TBs of data)

Fundamental unit of relational DBs: the record

Each entity in a DB has a corresponding **record**

- Features of a record are stored in **fields**
- Records with same “types” of fields collected into **tables**
- Each record is a row, each field is a column

ID	Name	UG University	Field	Birth Year	Age at Death
101010	Claude Shannon	University of Michigan	Electrical Engineering	1916	84
314159	Albert Einstein	ETH Zurich	Physics	1879	76
21451	Ronald Fisher	University of Cambridge	Statistics	1890	72

Table with six fields and three records.

Fields can contain different data types

ID	Name	UG University	Field	Birth Year	Age at Death
101010	Claude Shannon	University of Michigan	Electrical Engineering	1916	84
314159	Albert Einstein	ETH Zurich	Physics	1879	76
21451	Ronald Fisher	University of Cambridge	Statistics	1890	72

Unsigned int, String, String, String, Unsigned int, Unsigned int

Of course, can also contain floats, signed ints, etc. Some DB software allows categorical types (e.g., letter grades).

By convention, each record has a **primary key**

ID	Name	UG University	Field	Birth Year	Age at Death
101010	Claude Shannon	University of Michigan	Electrical Engineering	1916	84
314159	Albert Einstein	ETH Zurich	Physics	1879	76
21451	Ronald Fisher	University of Cambridge	Statistics	1890	72

Primary key used to uniquely identify the entity associated to a record, and facilitates joining information across tables.

ID	PhD Year	PhD University	Thesis Title
101010	1940	MIT	An Algebra for Theoretical Genetics
314159	1905	University of Zurich	A New Determination of Molecular Dimensions
21451			

ACID: Atomicity, Consistency, Isolation, Durability

Atomicity: to outside observer, every transaction (i.e., changing the database) should appear to have happened “instantaneously”.

Consistency: DB changes should leave the DB in a “valid state” (e.g., changes to one table that affect other tables are propagated before the next transaction)

Isolation: concurrent transactions don’t “step on each other’s toes”

Durability: changes to DB are permanent once they are committed

Note: some RDBMSs achieve faster performance, at cost of one or more of above

Related: Brewer’s Theorem https://en.wikipedia.org/wiki/CAP_theorem

Relational Database Management Systems (RDBMSs)

Program that facilitates interaction with database is called RDBMS

Public/Open-source options:

MySQL, PostgreSQL, **SQLite**

Proprietary:

IBM Db2, Oracle, SAP, SQL Server (Microsoft)

We'll use SQLite because it comes built-in to Python. More later.

SQL (originally SEQUEL, from IBM)

Structured Query Language (**Structured English QUery Language**)

Language for interacting with relational databases

Not the only way to do so, but by far most popular

Slight variation from platform to platform (“dialects of SQL”)

Good tutorials/textbooks:

https://www.w3schools.com/sql/sql_intro.asp

O'Reilly books: *Learning SQL* by Beaulieu

SQL Pocket Guide by Gennick

Severance, Chapter 14: <http://www.pythonlearn.com/html-270/book015.html>

Examples of database operations

ID	Name	GPA	Major	Birth Year
101010	Claude Shannon	3.1	Electrical Engineering	1916
500100	Eugene Wigner	3.2	Physics	1902
314159	Albert Einstein	4.0	Physics	1879
214518	Ronald Fisher	3.25	Statistics	1890
662607	Max Planck	2.9	Physics	1858
271828	Leonard Euler	3.9	Mathematics	1707
999999	Jerzy Neyman	3.5	Statistics	1894
112358	Ky Fan	3.55	Mathematics	1914

- Find names of all physics majors
- Compute average GPA of students born in the 19th century
- Find all students with GPA > 3.0

SQL allows us to easily specify queries like these (and far more complex ones).

Common database operations

Extracting records: find all rows in a table

Filtering records: retain only the records (rows) that match some criterion

Sorting records: reorder selected rows according to some field(s)

Adding/deleting records: insert new row(s) into a table or remove existing row(s)

Grouping records: gather rows according to some field

Adding/deleting tables: create new or delete existing tables

Merging tables: combine information from multiple tables into one table

Common database operations

SQL includes keywords for succinctly expressing all of these operations.

Extracting records: find all rows in a table

Filtering records: retain only the records (rows) that match some criterion

Sorting records: reorder selected rows according to some field(s)

Adding/deleting records: insert new row(s) into a table or remove existing row(s)

Grouping records: gather rows according to some field

Adding/deleting tables: create new or delete existing tables

Merging tables: combine information from multiple tables into one table

Retrieving records: SQL `SELECT` Statements

Basic form of a SQL `SELECT` statement:

```
SELECT [column names] FROM [table]
```

Example: we have table `t_customers` of customer IDs, names and companies

Retrieve all customer names: `SELECT name FROM t_customers`

Retrieve all company names: `SELECT company FROM t_customers`

Note: by convention (and good practice), one often names tables to be prefixed with “TB_” or “t_”. In our illustrative examples, I won’t always do this for the sake of space and brevity, but I highly recommend it in practice. See <https://launchbylunch.com/posts/2014/Feb/16/sql-naming-conventions/> and <http://leshazlewood.com/software-engineering/sql-style-guide/> for two people’s (differing) opinions.

Table t_students

id	name	gpa	major	birth_year	pets	favorite_color
101010	Claude Shannon	3.1	Electrical Engineering	1916	2	Blue
314159	Albert Einstein	4.0	Physics	1879	0	Green
999999	Jerzy Neyman	3.5	Statistics	1894	1	Red
112358	Ky Fan	3.55	Mathematics	1914	2	Green

```
SELECT id, name, birth_year FROM t_students
```

id	name	birth_year
101010	Claude Shannon	1916
314159	Albert Einstein	1879
999999	Jerzy Neyman	1894
112358	Ky Fan	1914

Filtering records: SQL WHERE Statements

To further filter the records returned by a `SELECT` statement:

```
SELECT [column names] FROM [table] WHERE [filter]
```

Example: table `t_inventory` of product IDs, unit cost, and number in stock

Retrieve IDs for all products with unit cost at least \$1:

```
SELECT id FROM t_inventory WHERE unit_cost >= 1
```

Note: Possible to do much more complicated filtering, e.g., regexes, set membership, etc. We'll discuss that more in a few slides.

Table t_students

id	name	gpa	major	birth_year	pets	favorite_color
101010	Claude Shannon	3.1	Electrical Engineering	1916	2	Blue
314159	Albert Einstein	4.0	Physics	1879	0	Green
999999	Jerzy Neyman	3.5	Statistics	1894	1	Red
112358	Ky Fan	3.55	Mathematics	1914	2	Green

```
SELECT id, name FROM t_students WHERE birth_year >1900
```

id	name
101010	Claude Shannon
112358	Ky Fan

NULL means Nothing!

Table t_thesis

id	phd_year	phd_university	thesis_title
101010	1940	MIT	An Algebra for Theoretical Genetics
314159	1905	University of Zurich	A New Determination of Molecular Dimensions
214511			
774477	1970	MIT	

```
SELECT id FROM t_thesis WHERE phd_year IS NULL
```

id
21451

NULL matches the *empty string*, i.e., matches the case where the field was left empty. Note that if the field contains, say, ` ` , then NULL will **not** match that row!

Ordering records: SQL ORDER BY Statements

To order the records returned by a `SELECT` statement:

```
SELECT [columns] FROM [table] ORDER BY [column] [ASC|DESC]
```

Example: table `t_inventory` of product IDs, unit cost, and number in stock

Retrieve IDs, # in stock, for all products, ordered by descending # in stock:

```
SELECT id, number_in_stock FROM t_inventory  
  
ORDER BY number_in_stock DESC
```

Note: most implementations order ascending by default, but best to always specify, for your sanity and that of your colleagues!

Table t_students

id	name	gpa	major	birth_year	pets	favorite_color
101010	Claude Shannon	3.1	Electrical Engineering	1916	2	Blue
314159	Albert Einstein	4.0	Physics	1879	0	Green
999999	Jerzy Neyman	3.5	Statistics	1894	1	Red
112358	Ky Fan	3.55	Mathematics	1914	2	Green

```
SELECT id, name, gpa FROM t_students ORDER BY gpa DESC
```

id	name	gpa
314159	Albert Einstein	4.0
112358	Ky Fan	3.55
999999	Jerzy Neyman	3.5
101010	Claude Shannon	3.1

More filtering: DISTINCT Keyword

To remove repeats from a set of returned results:

```
SELECT DISTINCT [columns] FROM [table]
```

Example: table `t_student` of student IDs, names, and majors

Retrieve all the majors:

```
SELECT DISTINCT major FROM t_student
```

Table t_students

id	name	gpa	major	birth_year	pets	favorite_color
101010	Claude Shannon	3.1	Electrical Engineering	1916	2	Blue
314159	Albert Einstein	4.0	Physics	1879	0	Green
999999	Jerzy Neyman	3.5	Statistics	1894	1	Red
112358	Ky Fan	3.55	Mathematics	1914	2	Green

```
SELECT DISTINCT pets FROM t_students ORDER BY pets ASC
```

Test your understanding: what should this return?

Table t_students

id	name	gpa	major	birth_year	pets	favorite_color
101010	Claude Shannon	3.1	Electrical Engineering	1916	2	Blue
314159	Albert Einstein	4.0	Physics	1879	0	Green
999999	Jerzy Neyman	3.5	Statistics	1894	1	Red
112358	Ky Fan	3.55	Mathematics	1914	2	Green

```
SELECT DISTINCT pets FROM t_students ORDER BY pets ASC
```

pets
0
1
2

More on WHERE Statements

WHERE keyword supports all the natural comparisons one would want to perform

(Numerical) Operation	Symbol/keyword
Equal	=
Not equal	<>
Less than	<
Less than or equal to	<=
Greater than	>
Greater than or equal to	>=
Within a range	BETWEEN ... AND ...

Examples:

```
SELECT id from t_student WHERE ...
```

```
... gpa>=3.2
```

```
... pets=1
```

```
... gpa BETWEEN 2.9 AND 3.1
```

```
... birth_year > 1900
```

```
... pets <> 0
```

Caution: different implementations define BETWEEN differently (i.e., inclusive vs exclusive)! Be sure to double check!

More on WHERE Statements

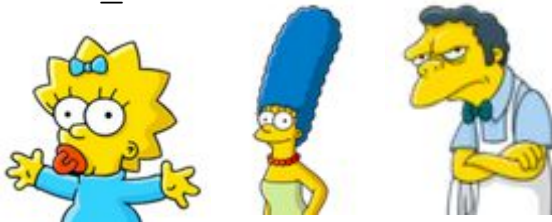
WHERE keyword also allows (limited) regex support and set membership

```
SELECT id, major from t_student WHERE major IN ("Mathematics","Statistics")
```

```
SELECT id, major from t_student WHERE major NOT IN ("Physics")
```

Regex-like matching with LIKE keyword, wildcards `_` and `\%`

```
SELECT id,name from t_simpsons_characters WHERE first_name LIKE "M%"
```



```
SELECT id,name from t_simpsons_characters WHERE first_name LIKE "B_rt"
```

Matches `'Bart'`, `'Bert'`, `'Bort'`...

Aggregating results: GROUP BY

I have a DB of transactions at my internet business, and I want to know how much each customer has spent in total.

customer_id	customer	order_id	dollar_amount
101	Amy	0023	25
200	Bob	0101	10
315	Cathy	0222	50
200	Bob	0120	12
310	Bob	0429	100
315	Cathy	0111	33
101	Amy	0033	25
315	Cathy	0504	70

```
SELECT customer_id, SUM(dollar_amount)
FROM t_transactions GROUP BY customer_id
```

customer_id	dollar_amount
101	50
200	22
310	100
315	153

GROUP BY `field_x` combines the rows with the same value in the field `field_x`

More about GROUP BY

GROUP BY supports other operations in addition to SUM:

COUNT, AVG, MIN, MAX

Called **aggregate** functions

Can filter results *after* GROUP BY using the HAVING keyword

```
SELECT customer_id, SUM(dollar_amount) AS total_dollar FROM t_transactions  
GROUP BY customer_id HAVING total_dollar > 50
```

customer_id	dollar_amount
101	50
200	22
310	100
315	153



customer_id	total_dollar
310	100
315	153

More about GROUP BY

GROUP BY supports other operations in addition to SUM:

COUNT, AVG, MIN, MAX

Called **aggregate** functions

Note: the difference between the HAVING keyword and the WHERE keyword is that HAVING operates *after* applying filters and GROUP BY.

Can filter results *after* GROUP BY using the HAVING keyword

```
SELECT customer_id, SUM(dollar_amount) AS total_dollar FROM t_transactions  
GROUP BY customer_id HAVING total_dollar > 50
```

customer_id	dollar_amount
101	50
200	22
310	100
315	153



customer_id	total_dollar
310	100
315	153

The AS keyword just lets us give a nicer name to the aggregated field.

Merging tables: JOIN

ID	Name	GPA	Major	Birth Year
101010	Claude Shannon	3.1	Electrical Engineering	1916
314159	Albert Einstein	4.0	Physics	1879
999999	Jerzy Neyman	3.5	Statistics	1894
112358	Ky Fan	3.55	Mathematics	1914

ID	#Pets	Favorite Color
101010	2	Blue
314159	0	Green
999999	1	Red
112358	2	Green

Join tables based on primary key

ID	Name	GPA	Major	Birth Year	#Pets	Favorite Color
101010	Claude Shannon	3.1	Electrical Engineering	1916	2	Blue
314159	Albert Einstein	4.0	Physics	1879	0	Green
999999	Jerzy Neyman	3.5	Statistics	1894	1	Red
112358	Ky Fan	3.55	Mathematics	1914	2	Green

Merging tables: INNER JOIN

t_student

id	name	gpa	major	birth_year
101010	Claude Shannon	3.1	Electrical Engineering	1916
314159	Albert Einstein	4.0	Physics	1879
999999	Jerzy Neyman	3.5	Statistics	1894
112358	Ky Fan	3.55	Mathematics	1914

t_personal

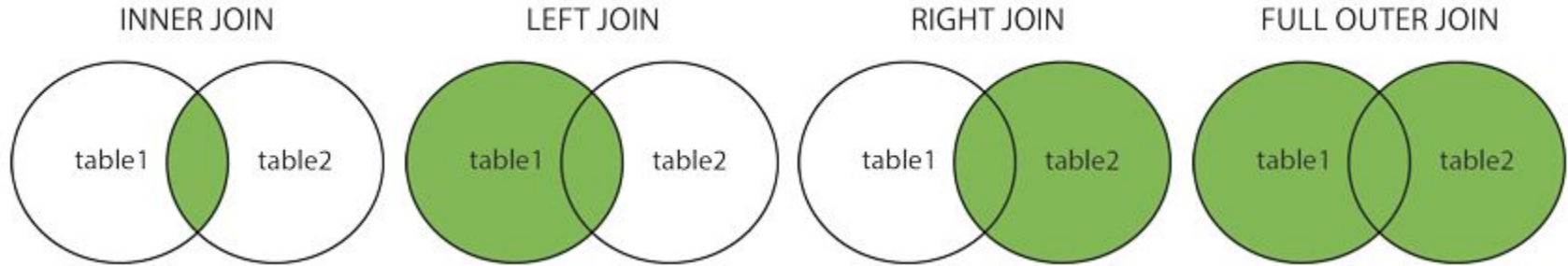
id	pets	favorite_color
101010	2	Blue
314159	0	Green
999999	1	Red
112358	2	Green

Join tables based on primary key

```
SELECT id, name, favorite_color
FROM
t_student INNER JOIN t_personal
ON t_student.id=t_personal.id
```

id	name	favorite_color
101010	Claude Shannon	Blue
314159	Albert Einstein	Green
999999	Jerzy Neyman	Red
112358	Ky Fan	Green

Other ways of joining tables: OUTER JOIN



(INNER) JOIN: Returns records that have matching values in both tables

LEFT (OUTER) JOIN: Return all records from the left table, and the matched records from the right table

RIGHT (OUTER) JOIN: Return all records from the right table, and the matched records from the left table

FULL (OUTER) JOIN: Return all records when there is a match in either left or right table

Creating/modifying/deleting rows

Insert a row into a table: `INSERT INTO`

```
INSERT INTO table_name [col1, col2, col3, ...]  
VALUES value1, value2, value3, ...
```


Note: if adding values for all columns, you only need to specify the values.

Modify a row in a table: `UPDATE`

```
UPDATE table_name SET col1=value1,col2=value2,  
WHERE condition
```

Delete rows from a table: `DELETE`

```
DELETE FROM table_name WHERE condition
```



Caution: if `WHERE` clause is left empty, you'll delete/modify the whole table!

Creating and deleting tables

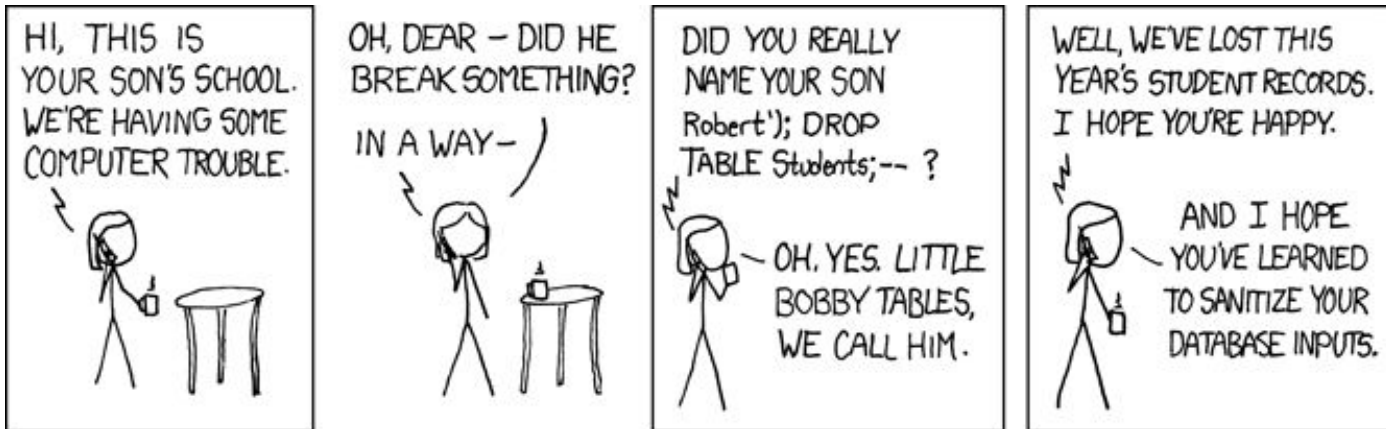
Create a new table: `CREATE TABLE`

```
CREATE TABLE table_name [col1 datatype, col2 datatype, ...]
```

Delete a table: `DROP TABLE`

```
DROP TABLE table_name;
```

Be careful when
dropping tables!



Python `sqlite3` package implements SQLite

`Connection` object represents a database

`Connection` object can be used to create a `Cursor` object

`Cursor` facilitates interaction with database

```
conn = sqlite3.connect('example.db')
```

establish connection to given DB file (creating it if necessary)

return `Connection` object

```
c = conn.cursor()
```

Creates and returns a `Cursor` object for interacting with DB

```
c.execute( [SQL command] )
```

runs the given command; cursor now contains query results

Python `sqlite3` package

Important point: unlike many other RDBMSs, SQLite does not allow multiple connections to the same database at the same time.

So, if you're working in a distributed environment, you'll need something else e.g., MySQL, Oracle, etc.

Python sqlite3 in action

```
1 import sqlite3
2 conn = sqlite3.connect('example.db')
3 c = conn.cursor() # create a cursor object.
4 c.execute('''CREATE TABLE t_student (id, name, field, birth_year)''')
5 students = [(101010, 'Claude Shannon', 'Electrical Engineering', 1916),
6             (500100, 'Eugene Wigner', 'Physics', 1902),
7             (314159, 'Albert Einstein', 'Physics', 1879),
8             (214518, 'Ronald Fisher', 'Statistics', 1890),
9             (662607, 'Max Planck', 'Physics', 1858),
10            (271828, 'Leonard Euler', 'Mathematics', 1707),
11            (999999, 'Jerzy Neyman', 'Statistics', 1894),
12            (112358, 'Ky Fan', 'Mathematics', 1914)]
13 c.executemany('INSERT INTO t_student VALUES (?, ?, ?, ?)', students)
14 conn.commit() # Write the changes back to example.db
15 for row in c.execute('''SELECT * from t_student'''):
16     print row
```

```
(101010, u'Claude Shannon', u'Electrical Engineering', 1916)
(500100, u'Eugene Wigner', u'Physics', 1902)
(314159, u'Albert Einstein', u'Physics', 1879)
(214518, u'Ronald Fisher', u'Statistics', 1890)
(662607, u'Max Planck', u'Physics', 1858)
(271828, u'Leonard Euler', u'Mathematics', 1707)
(999999, u'Jerzy Neyman', u'Statistics', 1894)
(112358, u'Ky Fan', u'Mathematics', 1914)
```

Python sqlite3 in action

```
1 import sqlite3
2 conn = sqlite3.connect('example.db')
3 c = conn.cursor() # create a cursor object.
4 c.execute('''CREATE TABLE t_student (id, name, field, birth_year)''')
5 students = [(101010, 'Claude Shannon', 'Electrical Engineering', 1916),
6             (500100, 'Eugene Wigner', 'Physics', 1902),
7             (314159, 'Albert Einstein', 'Physics', 1879),
8             (214518, 'Ronald Fisher', 'Statistics', 1890),
9             (662607, 'Max Planck', 'Physics', 1858),
10            (271828, 'Leonard Euler', 'Mathematics', 1707),
11            (999999, 'Jerzy Neyman', 'Statistics', 1894),
12            (112358, 'Ky Fan', 'Mathematics', 1914)]
13 c.executemany('INSERT INTO t_student VALUES (?, ?, ?, ?)', students)
14 conn.commit() # Write the changes back to example.db
15 for row in c.execute('''SELECT * from t_student'''):
16     print row
```

Create the table. Note that we need not specify a data type for each column. SQLite is flexible about this.

Insert rows in the table.

Note: sqlite3 has special syntax for parameter substitution in strings. Using the built-in Python string substitution is insecure-- vulnerable to SQL injection attack.

Executing a query returns an iterator over query results.

```
(101010, u'Claude Shannon', u'Electrical Engineering', 1916)
(500100, u'Eugene Wigner', u'Physics', 1902)
(314159, u'Albert Einstein', u'Physics', 1879)
(214518, u'Ronald Fisher', u'Statistics', 1890)
(662607, u'Max Planck', u'Physics', 1858)
(271828, u'Leonard Euler', u'Mathematics', 1707)
(999999, u'Jerzy Neyman', u'Statistics', 1894)
(112358, u'Ky Fan', u'Mathematics', 1914)
```


Python sqlite3 annotated

```
1 import sqlite3
2 conn = sqlite3.connect('example.db')
```

Establishes a connection to the database stored in `example.db`.

```
3 c = conn.cursor()
```

`cursor` object is how we interact with the database. Think of it kind of like the cursor for your mouse. It points to, for example, a table, row or query results in the database.

```
4 c.execute(''CREATE TABLE t_student (id, name, field, birth_year)'')
```

`cursor.execute` will run the specified SQL command on the database.

```
13 c.executemany('INSERT INTO t_student VALUES (?, ?, ?, ?)', students)
14 conn.commit() # Write the changes back to example.db
```

`executemany` runs a list of SQL commands.

`commit` writes changes back to the file. Without this, the next time you open `example.db`, the table `t_student` will be empty!

```
17 conn.close()
```


Close the connection to the database. Think of this like Python file `close`.

Metainformation: sqlite_master

Special table that holds information about the “real” tables in the database

```
1 import os, sqlite3
2 os.remove('example.db') #remove old version of the database.
3 conn = sqlite3.connect('example.db')
4 c = conn.cursor()
5 c.execute(''CREATE TABLE t_student (id, name, field, birth_year)'')
6 c.execute(''CREATE TABLE t_thesis (thesis_id, phd_title phd_year)'')
7 for r in c.execute(''SELECT * FROM sqlite_master''):
8     print r
```

Two tables, named
t_student and t_thesis



```
(u'table', u't_student', u't_student', 2, u'CREATE TABLE t_student (id, name, field, birth_year)')
(u'table', u't_thesis', u't_thesis', 3, u'CREATE TABLE t_thesis (thesis_id, phd_title phd_year)')
```

Retrieving column names in sqlite3

```
1 c.execute('''SELECT * from t_student''')
2 c.description
```



`description` attribute contains the column names; returned as a list of tuples for agreement with a different Python DB API.

```
(( 'id', None, None, None, None, None, None, None),
 ( 'name', None, None, None, None, None, None, None),
 ( 'field', None, None, None, None, None, None, None),
 ( 'birth_year', None, None, None, None, None, None, None))
```

```
1 [desc[0] for desc in c.description]
```

```
['id', 'name', 'field', 'birth_year']
```

Note: this is especially useful in tandem with the `mysql_master` table when exploring a new database, like in your homework!

Intermission

Next up...PyTorch

What is PyTorch?

It's a Python-based scientific computing package targeted at two sets of audiences:

1. A replacement for NumPy to use the power of GPUs
2. a deep learning research platform that provides maximum flexibility and speed

PyTorch Installation

PyTorch Build	Stable (1.4)			Preview (Nightly)	
Your OS	Linux		Mac		Windows
Package	Conda	Pip		LibTorch	Source
Language	Python 2.7	Python 3.5	Python 3.6	Python 3.7	C++ / Java
CUDA	9.2		10.1		None
Run this Command:	pip3 install torch torchvision				

See <https://pytorch.org/get-started/locally/>

Numpy

```
# -*- coding: utf-8 -*-
import numpy as np

# N is batch size; D_in is input dimension;
# H is hidden dimension; D_out is output dimension.
N, D_in, H, D_out = 64, 1000, 100, 10

# Create random input and output data
x = np.random.randn(N, D_in)
y = np.random.randn(N, D_out)

# Randomly initialize weights
w1 = np.random.randn(D_in, H)
w2 = np.random.randn(H, D_out)

learning_rate = 1e-6
for t in range(500):
    # Forward pass: compute predicted y
    h = x.dot(w1)
    h_relu = np.maximum(h, 0)
    y_pred = h_relu.dot(w2)

    # Compute and print loss
    loss = np.square(y_pred - y).sum()
    print(t, loss)

    # Backprop to compute gradients of w1 and w2 with respect to loss
    grad_y_pred = 2.0 * (y_pred - y)
    grad_w2 = h_relu.T.dot(grad_y_pred)
    grad_h_relu = grad_y_pred.dot(w2.T)
    grad_h = grad_h_relu.copy()
    grad_h[h < 0] = 0
    grad_w1 = x.T.dot(grad_h)

    # Update weights
    w1 -= learning_rate * grad_w1
    w2 -= learning_rate * grad_w2
```

PyTorch

```
import torch

dtype = torch.FloatTensor
# dtype = torch.cuda.FloatTensor # Uncomment this to run on GPU

# N is batch size; D_in is input dimension;
# H is hidden dimension; D_out is output dimension.
N, D_in, H, D_out = 64, 1000, 100, 10

# Create random input and output data
x = torch.randn(N, D_in).type(dtype)
y = torch.randn(N, D_out).type(dtype)

# Randomly initialize weights
w1 = torch.randn(D_in, H).type(dtype)
w2 = torch.randn(H, D_out).type(dtype)

learning_rate = 1e-6
for t in range(500):
    # Forward pass: compute predicted y
    h = x.mm(w1)
    h_relu = h.clamp(min=0)
    y_pred = h_relu.mm(w2)

    # Compute and print loss
    loss = (y_pred - y).pow(2).sum()
    print(t, loss)

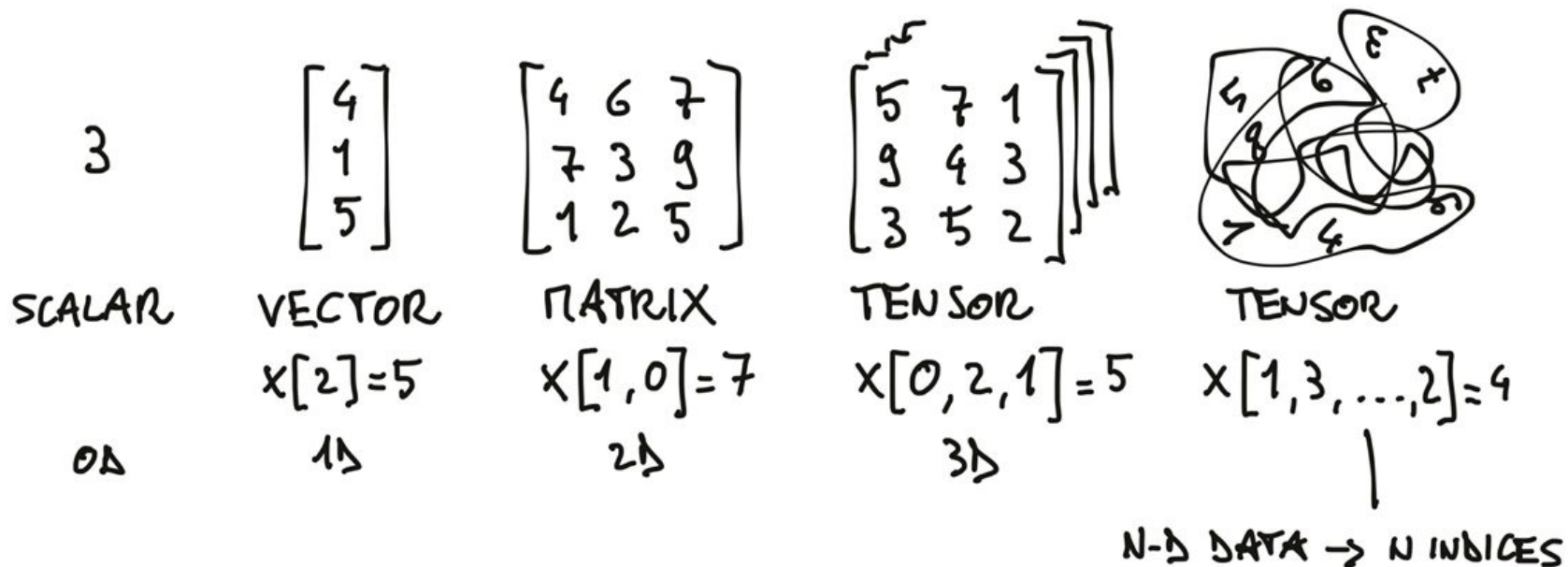
    # Backprop to compute gradients of w1 and w2 with respect to loss
    grad_y_pred = 2.0 * (y_pred - y)
    grad_w2 = h_relu.t().mm(grad_y_pred)
    grad_h_relu = grad_y_pred.mm(w2.t())
    grad_h = grad_h_relu.clone()
    grad_h[h < 0] = 0
    grad_w1 = x.t().mm(grad_h)

    # Update weights using gradient descent
    w1 -= learning_rate * grad_w1
    w2 -= learning_rate * grad_w2
```

numpy and
PyTorch have
many
similarities

Unit 1: PyTorch Tensors

Tensors generalize vectors/matrices, like NumPy's ndarrays



Credit: E. Stevens, L. Antiga, and T. Viehmann.
Deep Learning with PyTorch. 2020.

Properties of Tensors

Simply put, TENSORS are a generalization of vectors and matrices. In PyTorch, they are a multi-dimensional matrix containing elements of a **single data type**.

```
>>> import torch
>>> t = torch.tensor([[1., -1.], [1., -1.]])
>>> t
tensor([[ 1., -1.]
         [ 1., -1.]])

>>> t.dtype # They have a type
torch.float32

>>> t.shape # a shape
torch.Size([2, 2])

>>> t.device # and live in some device
device(type='cpu')
```

Creating Tensors

```
▶ import torch
```

Construct a 5x3 matrix, uninitialized:

```
[5] x = torch.empty(5, 3)  
    print(x)
```

```
tensor([[2.9901e-36, 0.0000e+00, 4.4842e-44],  
        [0.0000e+00,          nan, 0.0000e+00],  
        [6.5554e-10, 5.2479e-08, 4.2292e-05],  
        [1.7489e-04, 2.1858e-04, 3.3828e+21],  
        [1.3483e-05, 2.6728e+23, 5.0948e-14]])
```

Construct a randomly initialized matrix:

```
[ ] x = torch.rand(5, 3)  
    print(x)
```

Construct a matrix filled zeros and of dtype long:

```
[ ] x = torch.zeros(5, 3, dtype=torch.long)  
    print(x)
```

Construct a tensor directly from data:

```
[ ] x = torch.tensor([5.5, 3])  
    print(x)
```

Create a tensor based on an existing tensor

```
x = x.new_ones(5, 3, dtype=torch.double) # new_* methods take in sizes
print(x)

x = torch.randn_like(x, dtype=torch.float) # override dtype!
print(x) # result has the same size
```

Operations

Addition: syntax 1

 `x + y`

Addition: syntax 2

```
[ ] print(torch.add(x, y))
```

Addition: providing an output tensor as argument

```
[ ] result = torch.empty(5, 3)
    torch.add(x, y, out=result)
```

Addition: in-place

```
[ ] y.add_(x)
```

Resizing


Resizing: If you want to resize/reshape tensor, you can use `torch.view`:

```
[14] x = torch.randn(4, 4)
      y = x.view(16)
      z = x.view(-1, 8)  # the size -1 is inferred from other dimensions
      print(x.size(), y.size(), z.size())
```

 `torch.Size([4, 4]) torch.Size([16]) torch.Size([2, 8])`


If you have a one element tensor, use `.item()` to get the value as a Python number

```
▶ x = torch.randn(1)
  print(x)
  print(x.item())
```


 `tensor([-0.8768])`
`-0.8767884373664856`

Converting a Torch Tensor to and from a NumPy Array

```
[31] a = torch.ones(5)
      print(a)
```


 tensor([1., 1., 1., 1., 1.])

```
[32] b = a.numpy()
      print(b)
```


 [1. 1. 1. 1. 1.]

See how the numpy array changed in value.

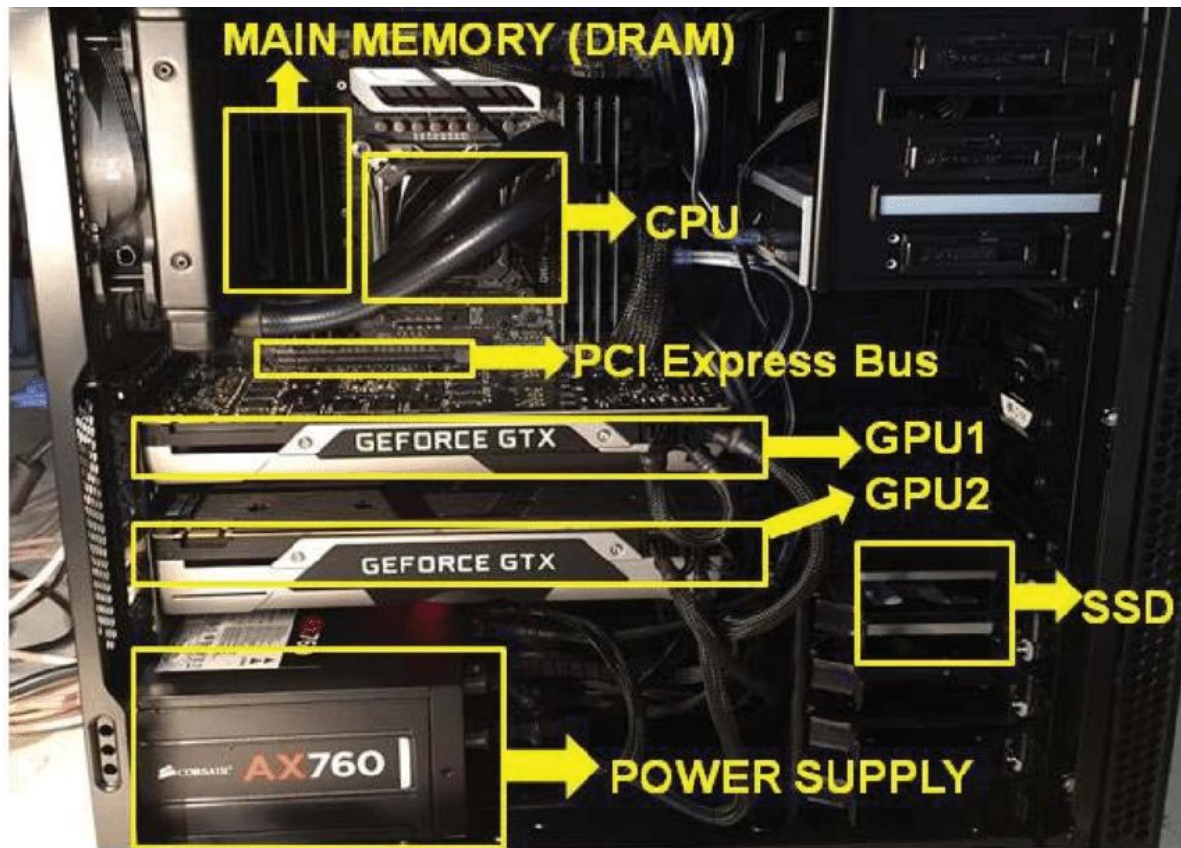
```
[33] a.add_(1)
      print(a)
      print(b)
```

 tensor([2., 2., 2., 2., 2.])
[2. 2. 2. 2. 2.]

```
[34] import numpy as np
      a = np.ones(5)
      b = torch.from_numpy(a)
      np.add(a, 1, out=a)
      print(a)
      print(b)
```

 [2. 2. 2. 2. 2.]
tensor([2., 2., 2., 2., 2.], dtype=torch.float64)

GPU computing



A GPU is connected to the motherboard. It is physical entity separate from the CPU.

Image credit: Tolga Soyata, from *GPU Parallel Program Development Using CUDA*

CUDA Tensors

```
[36] # let us run this cell only if CUDA is available
      # We will use ``torch.device`` objects to move tensors in and out of GPU
      if torch.cuda.is_available():
          device = torch.device("cuda")           # a CUDA device object
          y = torch.ones_like(x, device=device)   # directly create a tensor on GPU
          x = x.to(device)                        # or just use strings ``.to("cuda")``
          z = x + y
          print(z)
          print(z.to("cpu", torch.double))        # ``.to`` can also change dtype together!
```

x is now stored on the GPU

"cuda" stands for "Compute Unified Device Architecture"

GPU speedup on Google colab

```
[1] import torch  
import time
```

```
[2] x_cpu = torch.randn(120000, 10000)  
y_cpu = torch.randn(10000, 1)
```

```
[3] x_gpu = x_cpu.cuda()  
y_gpu = y_cpu.cuda()
```

```
[4] start = time.time()  
x_cpu.mm(y_cpu)  
print(time.time() - start)
```

0.39958739280700684

```
▶ start = time.time()  
x_gpu.mm(y_gpu)  
print(time.time() - start)
```

📄 0.004973173141479492

Matrix-vector multiplication on the GPU is nearly 100x faster!

Practice Problem 1:

Implementing an objective function in NumPy and PyTorch

Unit 2: Automatic Differentiation

Automatic Differentiation

```
[2] import torch
```

Create a tensor and set requires_grad=True to track computation with it

```
[3] x = torch.ones(2, 2, requires_grad=True)
    print(x)
```

```
tensor([[1., 1.],
        [1., 1.]], requires_grad=True)
```

Do an operation of tensor:

```
[4] y = x + 2
    print(y)
```

```
tensor([[3., 3.],
        [3., 3.]], grad_fn=<AddBackward0>)
```

y was created as a result of an operation, so it has a grad_fn.

```
[5] print(y.grad_fn)
```

```
<AddBackward0 object at 0x7f4f72142898>
```

Do more operations on y

```
[6] z = y * y * 3
    out = z.mean()
```

```
print(z, out)
```

```
tensor([[27., 27.],
        [27., 27.]], grad_fn=<MulBackward0>) tensor(27.,
```

Gradients

Let's backprop now Because `out` contains a single scalar, `out.backward()` is equivalent to `out.backward(torch.tensor(1))`.

```
[8] out.backward()
```

```
print gradients d(out)/dx
```

```
[9] print(x.grad)
```



```
tensor([[4.5000, 4.5000],  
        [4.5000, 4.5000]])
```

You should have got a matrix of 4.5. Let's call the `out Tensor` " o ". We have that $o = \frac{1}{4} \sum_i z_i$, $z_i = 3(x_i + 2)^2$ and $z_i|_{x_i=1} = 27$. Therefore,

$\frac{\partial o}{\partial x_i} = \frac{3}{2}(x_i + 2)$, hence $\frac{\partial o}{\partial x_i}|_{x_i=1} = \frac{9}{2} = 4.5$.

You can do many crazy things with autograd!

More gradients

```
[10] x = torch.randn(3, requires_grad=True)
```

```
    y = x * 2
    while y.data.norm() < 1000:
        y = y * 2

    print(y)
```

```
tensor([-1301.9290,   196.3312,  1150.0977], grad_fn=<MulBackward0>)
```

```
[11] gradients = torch.tensor([0.1, 1.0, 0.0001], dtype=torch.float)
      y.backward(gradients)
```

```
    print(x.grad)
```

```
tensor([1.0240e+02, 1.0240e+03, 1.0240e-01])
```

We don't always need gradients

You can also stop autograd from tracking history on Tensors with `.requires_grad=True` by wrapping the code block in `with torch.no_grad():`

```
[12] print(x.requires_grad)
      print((x ** 2).requires_grad)

      with torch.no_grad():
          print((x ** 2).requires_grad)
```



True
True
False

Practice Problem 2:

Computing gradients manually and automatically