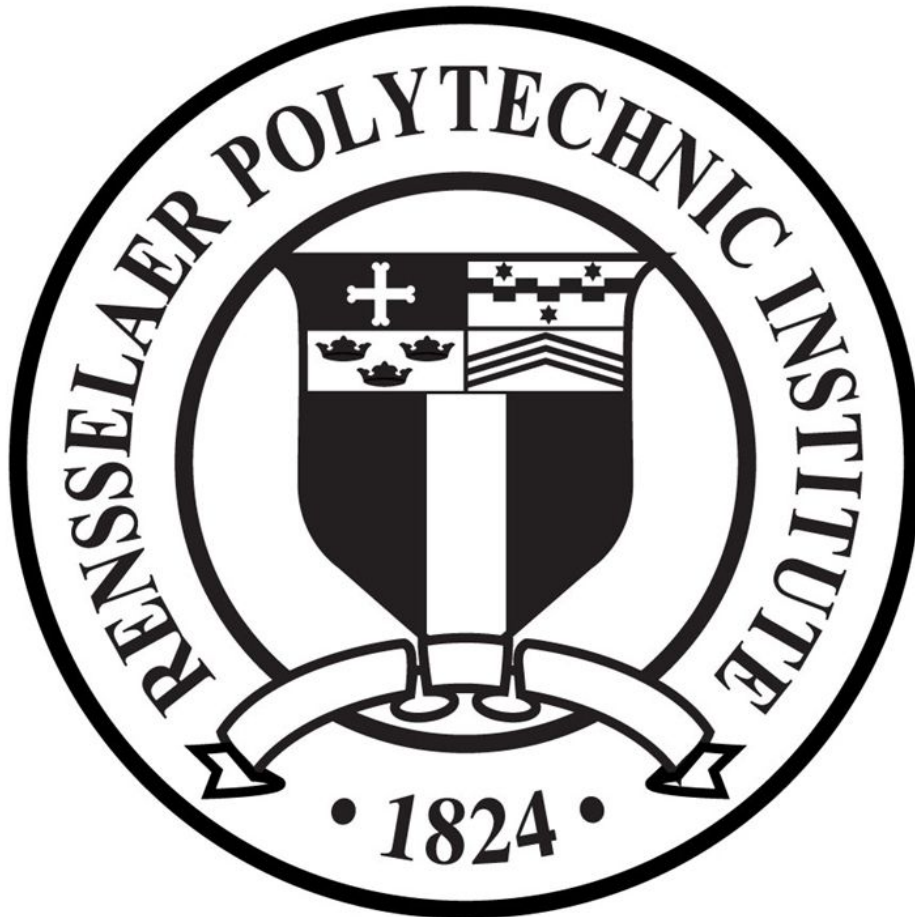


Project Schema

ITWS 6250 — Database Applications and Systems



Isaac Llewellyn & Shepard Gordon

Fall 2020

Summary

For this project, you will find multiple publicly available datasets that share common attributes (e.g., Zipcode), create a normalized schema describing the structure of the data, and produce an application that can populate your schema with the data—including the ability to refresh the data—and run queries on the data, producing useful output.

Objective

There are several objectives for this assignment.

- Gain an awareness of the scope of datasets publicly available for research purposes
- Demonstrate an ability to understand the structure of a dataset, as well as an ability to apply that understanding, using concepts learned during class, to create an effective database schema
- Apply concepts learned during class to query the data, and extend those concepts to create an application allowing users to do the same

Description

There are a number of different sources of publicly available data. Both the State of New York and the Federal Government provide hundreds of datasets. There are numerous other sources of open data as well, but those two will get you started. Please pay attention to licenses for any datasets you use. Data itself is generally not eligible for copyright protection (at least in the United States), but schemas are, and there may be terms of service for accessing the data itself.

Select two datasets that are robust enough to be interesting (a dataset with only four columns and a few thousand rows probably doesn't qualify). Students taking the course for graduate credit need to select two additional datasets, for a total of four.

They should share a common attribute (or set of attributes). Create a SQL schema for your data, making sure that it's appropriately normalized. Students taking the course for graduate credit need to use a non-relational database and schema for some portion of their data. Create an application in Python 3 that will load the dataset into a Postgres database defined by your schema. The loading process should be able to be re-run with updated datasets to refresh the data in the database.

Take some time to explore the data by running some SQL queries. Once you have an idea of some of the more interesting aspects of the data, create an interface for your application that will allow the user to explore the data as well.

Your application shouldn't re-implement the wheel. You don't need to provide the user with a way to do whatever they want. It should provide more of a self-guided tour, rather than a detailed map. It should provide interactivity beyond simply allowing the user to run one of five or six static queries, but it doesn't have to allow them to write their own queries.

For example, there might be a dataset giving the results of health inspections of restaurants in New York.

Your application might allow the user to see which restaurants in their area had violations, or how often a given restaurant received a violation, or whether restaurants in a certain area get more violations than other areas.

The interface can be text-based. If you want to go further and provide visualizations, that's fantastic, but it isn't within the scope of the project (you will not be graded on the appearance of your interface). Your application should be able to be built easily, the data loaded easily, and used easily.

You will demonstrate your application for the class in a short presentation in which you will discuss your choice of datasets, outline the design of your schema, and demonstrate the types of queries your application can perform.

All work will be done either individually or in teams of two.

Deliverables

There are four main deliverables:

- *Project Memo*
- *Database Schema*
- *Project Code*
- *Presentation*

Due Dates

- The memo is due on Submittity by 11:59pm on Friday October 9
- The schema is due on Submittity by 11:59pm on Wednesday October 26.
- The data-loading code is not formally due, but students should aim to have it completed by Wednesday November 25.
- The completed application is due on Submittity by 11:59 on Sunday December 6.
- You should be prepared to present your project to the class during the lecture period on Wednesday

Database Schema

You will submit a single SQL file that can be run to create the schema for your database. The SQL file will also be due before the rest of the project and will be graded at that time so that feedback can be incorporated into the final product.

Students taking the course for graduate credit will need to bear in mind that some portion of their data will be stored in a non-relational database.

That aspect of the project is not due with this deliverable.

Database Schema

- Team Members:

Isaac Llewellyn -- ITWS Cotermin Student
Shepard Gordon -- ITWS Graduate Student

- Datasets

[Current Season Spring Trout Stocking | State of New York](#)

Year INT	DEC Region INT	County TEXT	Town TEXT	Waterbody TEXT	Date TEXT	Number INT	Species Name TEXT	Size (inches) TEXT
-------------	-------------------	----------------	-----------	----------------	-----------	---------------	----------------------	------------------------

[National Register of Historic Places | State of New York](#)

Resource Name TEXT	County TEXT	National Register Date DATE	National Register Number TEXT	Longitude (number_???? Numeric_perh aps?????)	Latitude (number_??? ?Numeric_pe rhaps?????)	Location (location)
-----------------------	----------------	-----------------------------------	-------------------------------------	--	---	---------------------

[Recommended Fishing Rivers And Streams | State of New York](#)

Waterbody Name TEXT	Fish Species Present at Waterbody TEXT	Comments TEXT	Special Regulations on Waterbody TEXT	County TEXT	Types of Public Access s TEXT	Public Fishing Access Owner TEXT	Waterbody Information TEXT	Longitude (number_ ????Num eric_perh aps?????)	Latitude (number_ _????N umeric_ perhaps ?????)	Location (location)
---------------------------	---	------------------	---	----------------	--	--	----------------------------------	--	--	------------------------

[Fish Stocking Lists \(Actual\): Beginning 2011 | State of New York](#)

Year INT	County TEXT	Waterbody TEXT	Town TEXT	Month MONTH	Number (number_ ????Num eric_perh aps?????)	Species TEXT	Size (Inches) (number_????Numeric_p erhaps?????)
-------------	----------------	-------------------	--------------	----------------	---	-----------------	--

– How you plan to join the datasets

We plan on joining the datasets by their county fields, allowing people to explore the relationship between all New York fish stocking vs the subset of New York trout stocking, relating them in comparison to recommended historic places, fishing rivers and streams close to their location.

```
---
```

```
CREATE TABLE County_information (
```

```
County_name      TEXT ,  
Town_name        TEXT ,  
PRIMARY KEY (County_name, Town_name)
```

```
);
```

```
CREATE TABLE Stocking_information (
```

```
StockingID serial primary key,  
Year        INT,  
Waterbody   TEXT,  
Month        varchar(40),  
Number       INT,  
Species      TEXT,  
Size_Inches  TEXT,  
Future       boolean,  
County_name  TEXT,  
Town_name    TEXT,  
FOREIGN KEY (County_name, Town_name)  
REFERENCES County_information (County_name, Town_name)  
);
```

```
CREATE TABLE Waterbody_information (
```

```
Waterbody_Name      TEXT,  
Fish_Species_Present TEXT,  
Comments            TEXT,  
Special_Regulations TEXT,  
Types_of_Public_Access TEXT,  
Public_Fishing_Access_Owner TEXT,  
latitude            float,  
longitude           float,  
Location            POINT,  
Waterbody_information text,  
County_name         TEXT,  
PRIMARY KEY(Waterbody_Name, latitude, longitude)
```

```
);
```

```
CREATE TABLE County_historic (
```

```
Resource_Name      TEXT,  
National_Register_Date DATE,  
National_Register_Number TEXT PRIMARY KEY,  
Location           point,  
County_name        TEXT
```

```
);
```

```
--- lang=psql
```