# CSE 460 Sona Milestone 2

Jacob Goldverg

*Univeristy at Buffalo*

*UB# 50218299*

Team: Sona 460
Buffalo, US

jacobgol@buffalo.edu

*Abstract*—**The overview of the importer.py which is a python CLI that provides simple access to the John Hopkins Covid-19 dataset. The idea is to load many days into a single table, using a strict primary key.**

## I. PROBLEM STATEMENT

Johns Hopkins exposes daily covid reports on their GitHub, the issue we have found is that it is extremely difficult to manipulate this data inside as the files are all CSV's. The project will construct a database out of these files that will enable users to use this data in ways they were unable to in a simple Excel sheet. As there are new reports published for every country in the world every day it quickly becomes impossible to just use a simple Excel sheet to manage and manipulate these reports. Data surrounding COVID-19 varies heavily and changes very frequently which is why a database would be great to house the data, as well as manipulate it. This project will be a CLI application that provides simple arguments to build powerful queries that increase interpretability and manageability of the database.

## II. TARGET USER

The users of our database will be the general public, this means people who are already aware of how to connect into a database and manipulate some tables will be direct users. An example of a user is the frontend for our project as it will be able to do all of the operations specified above. Another example of an administrator is a Linux hobbyist that is trying to understand the data surrounding COVID-19. This would enable this person to study how covid affects their state to the entire United States.

The administration of the database will be done through a Python CLI application. The Docopt framework provides a simple CLI creation recipe and using Postgres SQL I was able to access the local deployment of Postgres and create a CLI application that makes it easy to insert/create/destroy/delete/update/alter operations on the database. The CLI is a completely public application on my Github account so anyone is able to use my code to learn more about the data. The CLI has checks to ensure that the operation cannot be harmful to the DB and the data within it. So, we are essentially configuring our Python CLI to be the sole administrator that many people can sign up to use and push all their requests through.

## A. Indicate the Primary key and the Foreign keys for each relation.

The database will have one tables as of right now, the core table is called "coviddays" and will house the daily reports provided by Johns Hopkins, and in this table the primary key is the UID of every row with the Last_Update timestamp. This way we can store multiple reports from the day in varying regions and not have to worry about a potential conflict between having many reports of the same day as a UID is completely unique. Any overlap in the primary key will result in that row being updated.

## B. Write the detailed description of each attribute(for each table), its purpose and datatype.

- Coviddays Table:
  - UID is an Integer, the unique row of the table
    - Cannot be null, part of primary key
  - Index: this is a new column in the table and represents the index of any given row and it is not the primary key
    - Cannot be null
  - Province_State varchar, the state or province the row is corresponding to.
    - Cannot be null and no default is available
  - Last_Update timestamp using UTF-8, the last time this measurement was updated
    - Cannot be null and no default is available and is part of the primary key
  - Lat: is latitude and is a Double Precision, initially we assumed we could use a Long.
    - Can be null and no default value
  - Long_: longitude: Double Precision, I initially assumed that I could use a Long.
    - Can be null and no default value
  - Confirmed: Integer, total number of cases this row(day).
    - Cannot be Null and has a default value of 0.
  - Deaths: Integer, total number of deaths due to covid.
    - Default is set to 0 and this cannot be null.
  - Recovered_Integer: total number of recovered cases this row(day).

- Default value is 0 and this value cannot be null.
  - o Active: Integer, total number of active cases this row(day).
    - Default value is 0 and this value cannot be null.
  - o Incident_Rate: Double Precision
    - Default value is 0 and this value cannot be null.
  - o Total_test_results Double Precision.
    - Default value is 0 and this value cannot be null.
  - o People_Hospitalized Double Precision.
    - Default value is 0 and this value cannot be null.
  - o Case_Fatality_Ratio Double Precision.
    - Default value is 0 and this value cannot be null.
  - o Testing_Rate Double Precision.
    - Default value is 0 and this value cannot be null.
  - o Hospitalization_Rate Double Precision.
    - Default value is 0 and this value cannot be null.
  - o ISO3: varchar that holds the three letters representing the ISO.
    - Cannot be null
  - o People_Hospitalized: A Double Precision value that holds the number of people hospitalized.
    - Can be null and has no default value.
  - o People_Tested: Double Precision value that holds the number of people tested on that day.
    - Cannot be null and has no default

The CLI application when importing csv files will automatically fill any null values with 0's, this way every single row will never have null values.

## III. INDICATE EACH ATTRIBUTES DEFAULT VALUE

This was done in step 2.

## IV. EXPLAIN THE ACTIONS TAKEN ON ANY FOREIGN KEY WHEN THE PRIMARY KEY IS DELETED.

For every entry we have a UID and Last_Updated date that allows every single row to be a totally unique row inside of our "coviddays" table. This way we can store many days regarding how states/provinces have been doing without worrying about uniqueness. The (UID, Last_Updated) is the primary key for the "coviddays" table. The Index column is a Secondary key which allows me to maintain the order of rows in the csv file as it is simply modular 50 for every entry in the file. Dropping any row from this table will have a cascading effect on the Views that the CLI creates. Provided below we have our "coviddays" table with a single day from the United States loaded into it. The Foreign key problem is not applicable to our schema design, the reason being is we have no foreign keys being used. But our schema does not support null and will "fix" every null to be 0, including for string data.

## V. THE CLI COMMANDS SUPPORTED

I decided that we only needed one table for this project called "coviddays". This table is responsible for housing all of the rows needed for the CLI to provide the user with clear and simple interface to manipulating and understanding the data from John Hopkins University.

The CLI focuses on providing a SQL free API that allows the user to still manage the database, and gain some statistics about the table. All the of SQL this project used can be found in the importer.py file where SQL statements are written into strings and then executed using a database connection. This is done by having a large number of commands where each command facilitates some new piece of functionality for the user.

Here is a list of all the commands the CLI supports with the corresponding functionality.

1. load_csv: Reads the csv file into a panda's DataFrame and then writes it to the coviddays table. Where if the table exists, it does a copy and concatenates the previous table to the csv data so it always preserves the data contained as well as add any extra columns if necessary.
2. get_date: Takes in three optional parameters (year, month, day) and returns all rows with the given combination of parameters.
3. drop_table: An administrator function that drops a table with the supplied name.
4. drop_column: Drops a column from the given table name as column name
5. alter_covid_table: This function checks that there is a primary key and unique constraint on the coviddays table. If there is no private key then it creates one using the columns (uid, last_updated) and creates a unique constraint on the exact same columns.
6. Stats: Takes in state, and column name and returns the average amount found in that column given the state
7. Print_table_stats: This takes in a table name with a default for coviddays. Ouputs: [(oid, schemaname, relname, heap_blks_read, heap_blks_hit, heap_blks_icache_hit, idx_blks_read, idx_blks_hit, idx_blks_icache_hit, toast_blks_read, toast_blks_hit, toast_blks_icache_hit, tidx_blks_read, tidx_blks_hit, tidx_blks_icache_hit] which are all different statistcis regarding the current state of the table. It also outputs each column, name, type and default value.
8. Sign_up: Takes in a username and password to create a postgres user that can access the coviddays table.
9. Sign_in: Takes in a username and password and uses them in the connection string to postgres. This ensures that every

user is authenticated with the DB before they can execute queries. If the user decides not to log in then the default user postgres will be used.

## VI. THE SQL STATEMENTS

To find the SQL statements used for any command please navigate to the GitHub repository and look at the importer.py file. There you will find functions with respective implementations.

To show some of the queries that I used with the respective functionality:

1. Create Primary Key Constraint: ALTER TABLE coviddays ADD CONSTRAINT pk_coviddays PRIMARY KEY (uid, last_update)
2. Create Unique constraint: ALTER TABLE coviddays ADD CONSTRAINT uni_coviddays UNIQUE(uid, last_update)
3. Sign up: CREATE USER name WITH PASSWORD pwd
4. Get rows on Date by day, month or year: SELECT * FROM coviddays WHERE EXTRACT(day FROM 'last_update') = day AND EXTRACT(month FROM 'last_update') = month AND EXTRACT(year FROM 'last_update') = year
5. How I check if a User already exists in the database: SELECT 1 FROM pg_roles WHERE rolname='username'

GitHub link: https://github.com/jgoldverg/dqmlProject1

## VII. LOADING LARGE DATA

The table is relatively trivial with only 17 dimensions but the large size comes from the number of day.csv files we can load into it. The queries done by the cli are not intensive as they do not apply large transformations.

Most of the operations are trivial in that they do simple create/update/delete/alter. This means that with a large number of days loaded into Postgres SQL the query time drops but by a logarithmic amount. The largest amount of files I added into the table is 30 which

## VIII. THE CONCLUSION

The project I created is for the use of managing the John Hopkins COVID-19 day to day files. This entailed creating a python script that handled all cli side which is the frontend of the project. It's a way for a user to directly get involved with my work. Then the CLI exposes various commands that execute different routines with arguments and flags. The cli follow POSIX style use, which means it is extremely standardized for arguments and flags.
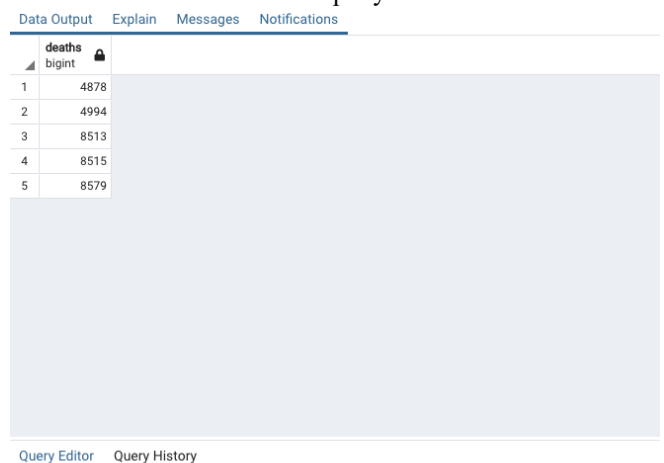
The database work for this assignment revolved around creating a key system that was simple and unique resistant this way not two rows from various kinds of files will collide. This allows the table to retain completely unique entries across any number of years while providing a

bucketing system using UID. The functionality my cli offers is provided above but does cover the general requirement of: create/load/alter/update/destroy. I also ensured security of the assignment by utilizing PostgresDB user/role system. The power of this, is every query ran by any user is logged so I am able to always find what queries users have done.

The future of this project would be to expand more on the kinds of statistics the cli is capable of showing. It would be nice to expand it to show maybe the std deviation and other kinds of metrics as data gets put in. With the goal to continue simplifying the way this data is used in a transparent way.

## IX. THE IMAGES

1. Here we have an image that shows the output from all deaths in Alabama for any day. We can tell that because we have 5 rows then are 5 days in the database from when this query was ran.



2. Here we have the E/R Diagram used for our coviddays table.

3. This image the database with a single day of covid testing results.



4. The row shere show a simple query that gets all the values with a total_test_results > 0.