

Installing the program:

Downloadable zip file can be found from: <https://github.com/PetKosTuni/Epidemic-model-of-COVID/releases> after the release of this program. Depending on errors either use python3 or python.

Windows and PowerShell (Python3 and pip already installed):

```
git clone git@github.com:PetKosTuni/Epidemic-model-of-COVID.git
```

(https url: <https://github.com/PetKosTuni/Epidemic-model-of-COVID.git>)

```
cd Epidemic-model-of-COVID
```

```
pip install virtualenv
```

```
python -m venv venv
```

```
.\venv\Scripts\activate
```

```
cd code
```

```
pip install -r requirements.txt
```

How to operate the program, program flow:

The user has to first move to the code directory: `cd code`

The command line would look like this ...Epidemic-model-of-COVID\code>

The program works by running a command in the command line (examples below) for first `validation.py`, and after it has finished, for `generate_predictions.py`.

The format for the command is 'python validation.py/generate_predictions.py (inputs)'. The inputs are first selected from the changeable inputs list: `--{input} {wanted input}`. **MAKE SURE THAT THE INPUTS ARE VALID.** For example, don't put end date before start date. The program will crash otherwise. Data has to exist between start date and mid date.

When generating predictions, the user should only switch the script filename from `validation.py` to `generate_predictions.py` keeping the same input arguments. The inputs for both `validation.py` and `generate_predictions.py` must be the same.

When running the files, csv-results are generated in directories beginning with *val_results* and *pred_results*. Plots are generated in directories beginning with *figure* in a pdf form.

Example commands for generating predictions:

Example 1:

Create predictions for all states using NYtimes dataset with end date 2021-07-07 and validation end date 2021-07-14.

Step 1: `python validation.py --END_DATE 2021-07-07 --VAL_END_DATE 2021-07-14 --dataset NYtimes --level state`

Step 2: `python generate_predictions.py --END_DATE 2021-07-07 --VAL_END_DATE 2021-07-14 --dataset NYtimes --level state`

Example 2:

Create predictions for Orange county in California using a custom dataset located in file path 'data/custom_dataset.csv'. Same dates are used as in example 1.

Step 1: `python validation.py --END_DATE 2021-07-07 --VAL_END_DATE 2021-07-14 --dataset CUSTOM_DATASET --dataset_filepath 'data/custom_dataset.csv' --level county --state California --county Orange`

Step 2: `python generate_predictions.py --END_DATE 2021-07-07 --VAL_END_DATE 2021-07-14 --dataset CUSTOM_DATASET --dataset_filepath 'data/custom_dataset.csv' --level county --state California --county Orange`

Examples with pictures on page 5.

Changeable inputs:

All dates in format yyyy-mm-dd. The inputs are the same for both files. Values must be given for END_DATE and VAL_END_DATE as other arguments can use default values.

--START_DATE

The date when the learning for the prediction starts. As default the starting date is taken from prediction_data.py file. **Some single county validations dont work without giving an early start date**, since the start date gets set later than the get-function end date (due to the counties not having enough data?).

--MID_DATE

The date when the second learning period starts. Must be at least 3 weeks before END_DATE and must be set by user if START_DATE is set by user. As 'default' the date is taken from prediction_data.py file.

--RESURGE_DATE

A date used when going through US regions. As 'default' the date is taken from prediction_data.py file. Should be between MID_DATE and END_DATE.

--END_DATE

The date when the final learning period ends. Must always be set by the user.

--VAL_END_DATE

The date when the prediction starts. Must always be set by the user. Should also be at least a week after END_DATE

--level

The level on which the program operates (nation/state/county). By default, the parameter is state.

--state

Validation/prediction for one specific state (level should be set as state).

--nation

Validation/prediction for one specific nation (level should be set as nation).

--county

Validation/prediction for one specific county (level should be set as county and county --state should also be specified)

--dataset

The user can choose which dataset to use. default = "NYtimes", if CUSTOM_DATASET is chosen, the user has to also input the filepath of the dataset (csv). Otherwise, if not NYtimes or CUSTOM_dataset. the dataset JHU is chosen.

--dataset_filepath

If a custom dataset has been chosen, by inputting dataset as CUSTOM_DATASET, the filepath of the dataset csv-file has to be put in the form 'data/...'

--pop_in

The amount of people joining the susceptible population. As default determined by the region between 0.0002 and 0.1.

--bias

The bias of the learning model. As default determined by the region studied between 0.1 and 0.5.

--pred_range

Range for prediction in days (default: 100 days).

If wanted, you can add your own changeable inputs by adding them in to the args-list in start of the two main files (validation/generate_predictions).

Prediction_data.py

The file contains many hardcoded values that can be changed if wanted, such as the number of valid countries that can be included in a custom dataset.

The prediction_data.py contains hard-coded values that cause restrictions / have to be changed if new rules are wanted for datasets! For example, currently 26 countries MUST BE IN THE custom dataset, since the code checks the country names from this file. You probably should add wanted countries to this file if they are included in a custom dataset, or remove them, if they are not included in the dataset. There may be more of these kind of restrictions, so reading the code thoroughly may be beneficial. For example, the US county Lassen does not work with the code for some reason (relic from earlier developers).

Custom datasets:

The datasets should be in csv-form and added to the data/ directory.

If a new dataset is wanted to be implemented as an csv, the **The DATASET_template should be able to create useful training data, if it includes confirmed cases, deaths** (and recovery amounts if wanted, but we don't think these are used anywhere). **The names for these columns may be slightly different in different dataset.csv-files**, so they can be changed in either the csv-files or by changing the hard coded column names in the prediction_data.py-file. This template dataset class should fit the generic use case for a dataset, where the dataset has the following (or similar columns):

```
custom_dataset_columns = ['date', 'country', 'state', 'county', 'cases', 'deaths', 'recoveries']
```

```
date, county, fips, cases, deaths
2020-01-21, Orange, 53, 1, 0
```

The dataset will pick a column and the values corresponding it. HOWEVER, **the template datasets have different restrictions depending on the level used.** For example, if the level is state, and there is county data, there would be multiple of the same state for one date. This would crash the program (perhaps because the array is simply too big in relation to some expected value?). In short, **with the DATASET_template, for state and nation levels, use one date per state/nation. County data will also require an additional state column which to match the county to.** It would also be wise to separate county/state/nation data in to their own csv-files, even though the template tries to check the correct level from the file. If using a single region to validate, these restrictions may not be so strict.

Try to also avoid false names for regions, since the populations are collected from certain csv-files, which have hard coded state, county and nation population data.

You can also add your own datasets as their own class with their own functions in the data.py, by using the other datasets as example, if the DATASET_template is not working satisfactorily, or you want to add other than csv-files with other kinds of columns. You also would have to add the datasets to the validation/generate_predictions-files as options in the data = ... rows (like where DATASET_template is used):

```
if args.dataset == "NYtimes":
    data = NYTimes(level='counties')
elif args.dataset == "CUSTOM_DATASET":
    data = DATASET_template(args.dataset_filepath, pdata.custom_dataset_columns, level = 'counties')
else:
    data = JHU_US(level='counties')
```

Examples with pictures:

Example 3:

Create predictions for Argentina (using JHU dataset as default) with end date 2021-07-07 and validation end date 2021-07-14. Prediction range is set as 300 days.

Running the following commands results in the graph below.

Step 1: cd code

```
PS T:\BACKUP\Koulu\Epidemic-model-of-COVID> cd code
PS T:\BACKUP\Koulu\Epidemic-model-of-COVID\code> █
```

Step 2: python validation.py --START_DATE 2020-04-06 --MID_DATE 2021-05-07 --END_DATE 2021-07-07 --VAL_END_DATE 2021-07-14 --level nation --nation Argentina --pred_range 300

```
PS T:\BACKUP\Koulu\Epidemic-model-of-COVID\code> python validation.py --START_DATE
2020-04-06 --MID_DATE 2021-05-07 --END_DATE 2021-07-07 --VAL_END_DATE 2021-07-14
--level nation --nation Argentina --pred_range 300 █
```

Validation results are saved in code/val_results_world directory.

▼ val_results_world

```
testJHU_val_params_best_END_DATE_2021-07-07_VAL_END_DATE_2021-07-14 U
testJHU_val_params_END_DATE_2021-07-07_VAL_END_DATE_2021-07-14 U
```

Validation result files are shown below.

```
code > val_results_world > testIHU_global_val_params_best_END_DATE_2021-07-07_VAL_END_DATE_2021-07-14
1 {"Argentina": [9039154.8, 11298.943500000001, 1000000.0, 6612837.4548219485, 114241.0113141035, 18321.25814324245, 0.3326389573748209, 0.0002471586336281687]}
```

```
code > val_results_world > testIHU_global_val_params_END_DATE_2021-07-07_VAL_END_DATE_2021-07-14
1 [1003218577435], [9039154.8, 11298.943500000001, 1000000.0, 6612837.4548219485, 114241.0113141035, 18321.25814324245, 0.3326389573748209, 0.0002471586336281687], [9039154.8, 11298.943500000001, 1000000.0, 6612837.4548219485, 114241.0113141035, 18321.25814324245, 0.3326389573748209, 0.0002471586336281687]
```

Step 3: python generate_predictions.py --START_DATE 2020-04-06 --MID_DATE 2021-05-07 --END_DATE 2021-07-07 --VAL_END_DATE 2021-07-14 --level nation --nation Argentina --pred_range 300

```
PS T:\BACKUP\Koulu\Epidemic-model-of-COVID\code> python generate_predictions.py
--START_DATE 2020-04-06 --MID_DATE 2021-05-07 --END_DATE 2021-07-07 --VAL_END_DATE 2021-07-14 --level nation --nation Argentina --pred_range 300
```

Prediction results are saved in code/pred_results_world directory.

The plots generated are saved in code/figure_nation directory.

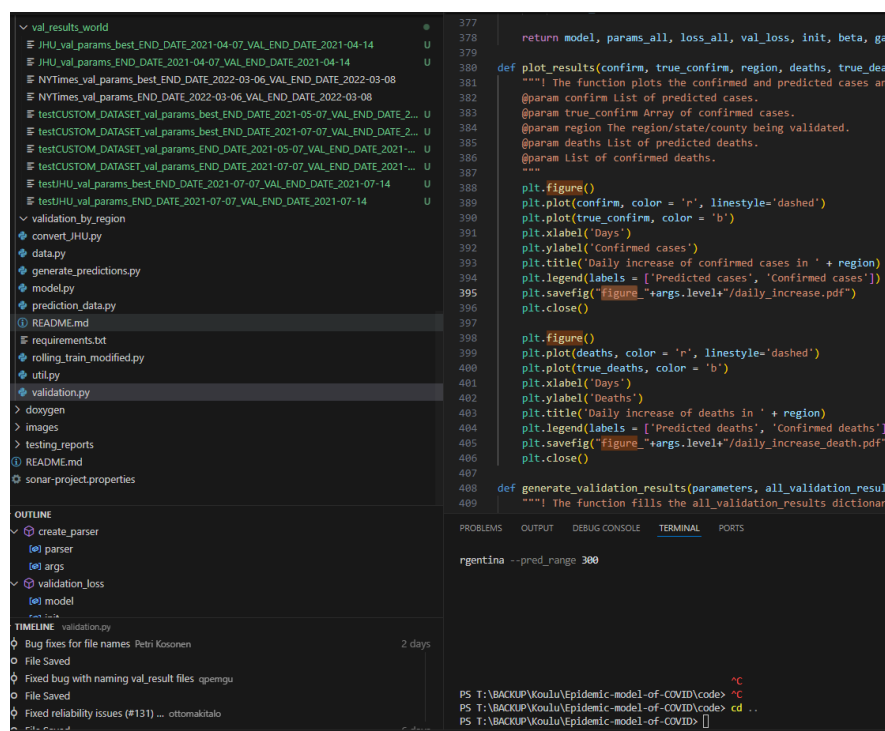
✓ pred_results_world

pred_nation_END_DATE_2022-03-06_PRED_START_DATE_2022-03-08.csv

```
code > pred_results_world > pred_nation_END_DATE_2021-07-07_PRED_START_DATE_2021-07-14.csv > data
1 Date,lower_pre_confirm,upper_pre_confirm,lower_pre_act,upper_pre_act,lower_pre_confirm_daily,lower_pre_act_daily,upper_pre_confirm_daily,upper_pre_act_daily,lower_pre_confirm_daily,lower_pre_act_daily,upper_pre_confirm_daily,upper_pre_act_daily
2 2021-07-14,4392657.0,4702657.0,100350.0,100350.0,4702657.0,100350.0,0.0,0.0,4702657.0,100350.0,0.0,0.0,Argentina
3 2021-07-15,4718213.285714285,4718213.285714285,100651.57142857143,100651.57142857143,4718213.285714285,100651.57142857143,401.5714285714348,401.5714285714348,15554.449711158313,288968.7954675979,36464.546393518
4 2021-07-16,4725008.322407215,4735579.388681255,101027.29393739869,101055.21310436636,4729455.678969429,101033.65127277929,382.07984426785215,375.7225088272535,403.6417657949286,24495.545174453237,284658.943908307,57113.85767951
```

✓ figure_nation
daily_increase_Argentina.pdf

This example will be thorough. First, we start in our terminal of choice once the program has been installed with the instructions from the start of the file. In Visual Studio Code the view should look like this after opening the program folder. The terminal is in the lower right corner.



First we need to move to the correct directory, code.

```
PS T:\BACKUP\Koulu\Epidemic-model-of-COVID> cd code
PS T:\BACKUP\Koulu\Epidemic-model-of-COVID\code> |
```

Before we can get the prediction results for the date range and region, we need to generate validation files. We have to decide the dataset, the dates for which the validation occurs, and the level of validation (county/state/nation). For this case we choose California, US on a state level with the END_DATE 2021-07-07 and VAL_END_DATE 2021-07-14. Rest of the inputs will be default values. The dataset will default to NYtimes. This will be a default run of the program, with only the mandatory inputs. If using the CUSTOM_dataset, remember to add the filepath and remember the restrictions written earlier.

The command to generate the wanted validation files would look like this (single line will do):

```
PS T:\BACKUP\Koulu\Epidemic-model-of-COVID\code> python validation.py
--END_DATE 2021-07-07 --VAL_END_DATE 2021-07-14 --state California |
```

The terminal will start printing information. Once the total confirmed cases and deaths are printed, validation files have been created.

```
Total confirmed cases: 4147599.9024823736 Total deaths: 67927.2987791542
PS T:\BACKUP\Koulu\Epidemic-model-of-COVID\code> |
```

Since our validation level was at state level, the validation files were written into /val_results_state, in csv-form:

```

  ✓ val_results_state
  └─ testNYtimes_val_params_best_END_DATE_2021-07-07_VAL_END_DATE_2021-07-14
  └─ testNYtimes_val_params_END_DATE_2021-07-07_VAL_END_DATE_2021-07-14

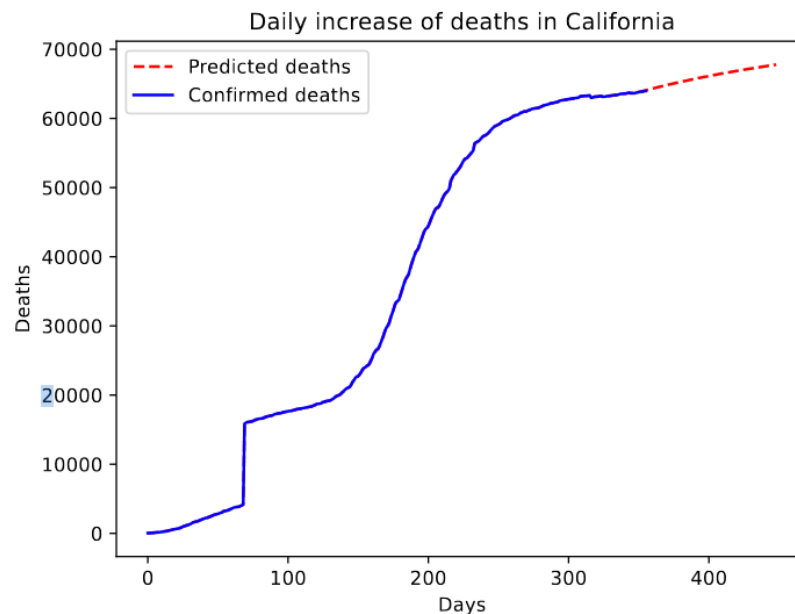
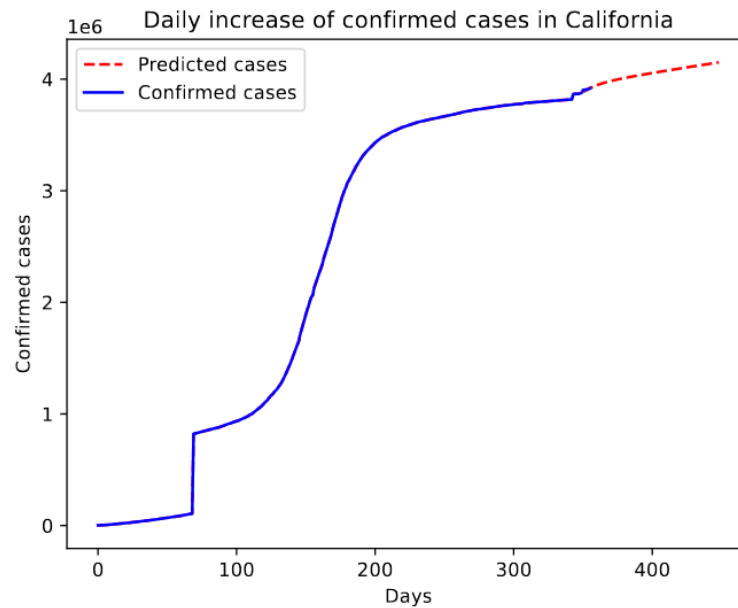
code > val_results_state > testNYtimes_val_params_best_END_DATE_2021-07-07_VAL_END_DATE_2021-07-14
1  {"California": [7902444.600000001, 112892.06571428572, 2.5608062424046038e-05, 4147599.9024823736, 67927.2987791542, 8393.714285714086

code > val_results_state > testNYtimes_val_params_END_DATE_2021-07-07_VAL_END_DATE_2021-07-14
1  {"California": [[7902444.600000001, 263414.82, 2.822198831836456e-05, 4135289.493174593, 67781.93515160913, 8393.714285714086, 0.00947
```

The validation process will also plot the latest region validated in the correct /figure_... folder, in our case in /figure_state. The plots are always called daily_increase_death.pdf and daily_increase.pdf:

```

  ✓ figure_state
  └─ daily_increase_Connecticut
  └─ daily_increase_death.pdf
  └─ daily_increase.pdf
```

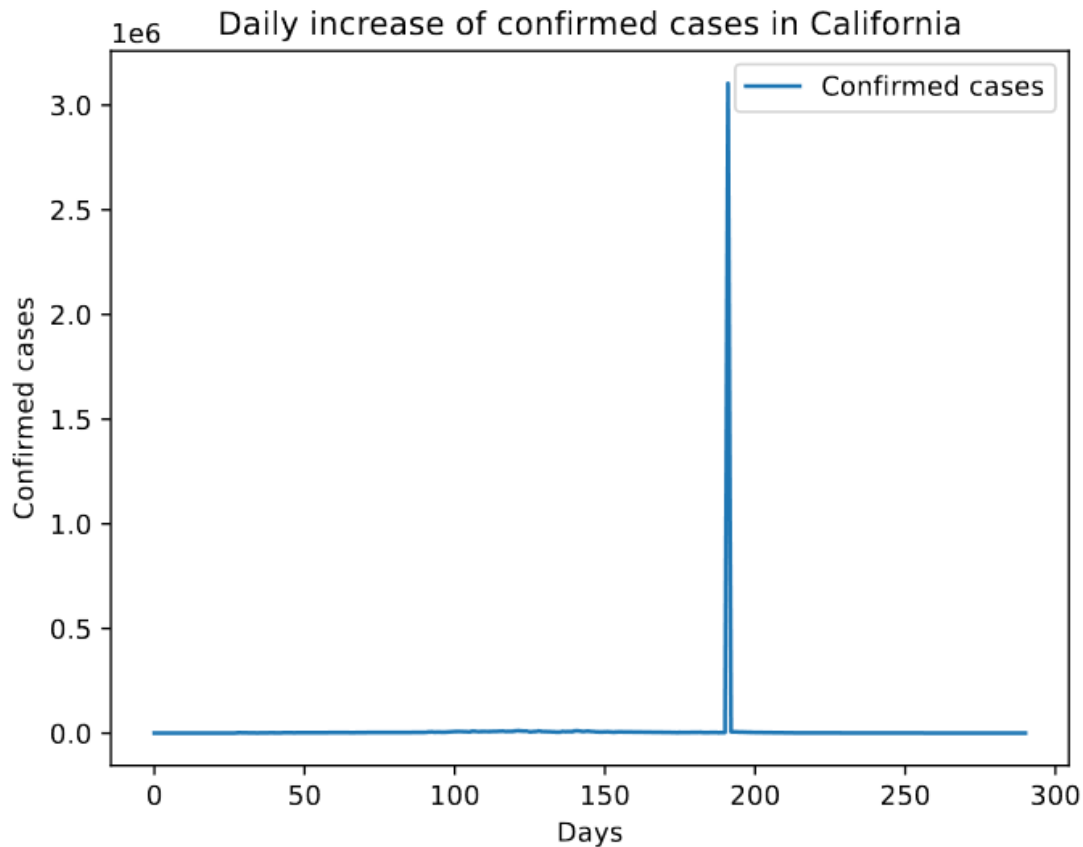
Now we can get the wanted prediction results. The command should be the exact same, expect for the file name. The command to generate the wanted prediction files looks like this (single line will do):

```
PS T:\BACKUP\Koulu\Epidemic-model-of-COVID\code> python generate_predictions.py
--END_DATE 2021-07-07 --VAL_END_DATE 2021-07-14 --state California
```

Once the terminal has printed the following information

```
region: California training loss: [0.009866025960771942, 0.00079612604486260
47, 0.033021082650520145] [0.006372791178703813, 0.02882654860777848] maximum
death cases: 69203 maximum confirmed cases: 4056809
```

the prediction results are generated. The program has plotted the Californian predictions in /figure_state, with the name daily_increase_California.pdf. The plot has a spike (see notes section):



The actual results are saved in csv-form in `> pred_results_state`:

```

pred_state_END_DATE_2021-07-07_PRED_START_DATE_2021-07-14.csv

code > pred_results_state > pred_state_END_DATE_2021-07-07_PRED_START_DATE_2021-07-14.csv > data
1 Date,lower_pre_confirm,upper_pre_confirm,lower_pre_fata,upper_pre_fata,pre_confirm,pre_fata,pre_fata_daily,lowe
18 2021-07-30,3982430.49261387,3998656.68617102,64616.19563550431,64952.55385520641,3984615.168763842,64832.073570
19 2021-07-31,3984338.93702625,4002106.0518392557,64649.3483090278,65023.455571137725,3986745.75481126,64888.94062
20 2021-08-01,3986133.8206677847,4005452.031852577,64681.81627992996,65094.97286795331,3988761.605106856,64945.875
21 2021-08-02,3987825.0748998797,4008699.8576953514,64713.63468358342,65167.06029947563,3990681.406205948,65002.86
22 2021-08-03,3989421.7192702545,4011854.3760434124,64744.83684502908,65239.67470580893,3992509.9080927237,65059.7
23 2021-08-04,3990931.947301081,4014920.0830881777,64775.45436604252,65312.77509742698,3994249.237580806,65116.620
24 2021-08-05,3992363.204046904,4017901.155528006,64805.51721141901,65386.32254880214,3995919.5261704447,65173.429
25 2021-08-06,3993722.256257848,4020801.4785962394,64835.05379177596,65460.28009804344,3997459.5979060354,65230.19
26 2021-08-07,3995015.2558632395,4023624.6714193416,64864.091042452696,65534.61265159866,3998959.4988414547,65286.
27 2021-08-08,3996247.7974135126,4026374.1099614454,64892.654498582546,65609.28689363552,4000437.937998033,65343.5
28 2021-08-09,3997424.970055713,4029052.9477853607,64920.75550403625,65684.27119985354,4001862.832163989,65400.097
29 2021-08-10,3998551.404562884,4031664.134837905,64948.41113467042,65759.53555550994,4003238.542526752,65456.5576
30 2021-08-11,3999631.3158880034,4034210.434447572,64975.65679784629,65835.05147744711,4004542.0154162105,65512.99
31 2021-08-12,4000668.6157607487,4036694.438704659,65002.51286251357,65910.79193993488,4005783.0888361703,65569.56
32 2021-08-13,4001666.815459294,4039118.582377778,65028.99862778827,65986.73130413175,4006972.7132765236,65626.031
33 2021-08-14,4002628.9735086183,4041485.155506041,65055.13237852696,66062.84442218767,4008142.521874979,65682.403
34 2021-08-15,4003557.874693424,4043796.3147929367,65080.93143818704,66139.10593975116,4009256.084313233,65738.668
35 2021-08-16,4004456.05646849,4046054.114458737,65106.41221910132,66215.4942936489,4010376.6254979335,65794.82076
36 2021-08-17,4005325.8323127767,4048260.7387470235,65131.500270303685,66291.9892152242,4011477.800350679,65850.85

```

Notes:

rolling_train_modified.py has some unreachable code in the loss train function for unknown reason. You can try fixing the code if it makes sense and makes better predictions.

validation.py and generation_predictions.py have resurge_date if-statements, that cause spikes in the graph results. The if statements and the resurge dates could just be removed to fix this problem. Currently this problem only exists for US counties/states.

If removed files or functionalities are required, the old code from the original project is still available.

The tutorial and video demo will be included in the handoff email and the root path of the repository. The /code-folder will include another read-me that also includes some of the information from the tutorial.

The test folder includes a readme to run the unit tests.

The doxygen folder includes a file called index.html, which when opened in a browser, includes a useful list of classes used in the program.