

Marc Laugharn

Covid-19 Lockdown Impact Model

I was interested in finding out the impact of lockdown on the acceleration of the daily cases per state in the United States (and Washington DC). I collected timeseries data from a wide variety of sources:

- Apple's COVID-19 Mobility Trends Data, specifically change in percent of driving from baseline, as recorded by Apple Maps usage
- Google's COVID-19 Mobility Trends Data, as percent change from baseline:
 - Retail and recreation
 - Grocery and pharmacies
 - Parks
 - Transit stations
 - Workplaces
 - Residential areas
- COVID Act Now's daily data, including:
 - Hospital beds required and in use
 - ICU Beds capacity and in use
 - Ventilators capacity and in use
 - Real time estimate of R_0 value, R_t (unsure how this value is computed.. I asked Covid19ActNow but they were ambiguous and deflected to look at an upcoming blog post)
 - Cumulative dead, infected
 - Cumulative positive, negative tests
- Lockdown level by state, where:
 0. No or few containment measures
 1. Ban on public gatherings, cancellation of major events
 2. Schools and universities closed
 3. Nonessential shops, restaurants and bars closed
 4. Night curfew/partial lockdown
 5. All-day lockdown: shelter in place order, citizens allowed to leave home
 6. Harsh lockdown: citizens not allowed to leave home

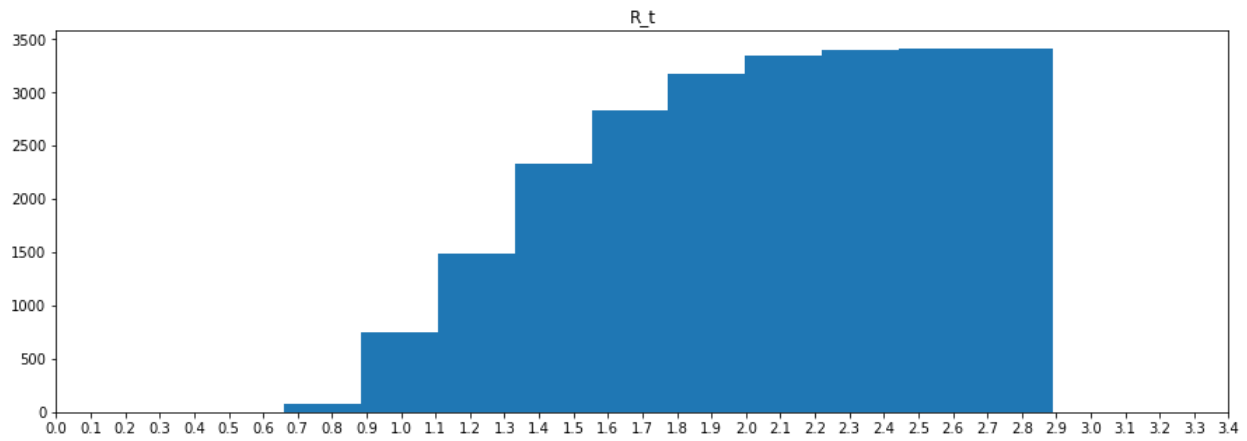
I also collected various non-time-series data per region:

- US Census health and poverty data (including obesity rate, smoking rate, income, poverty rate, healthcare access)
- Density per square mile
- 2018 population
- Percentage of people living in urban environments

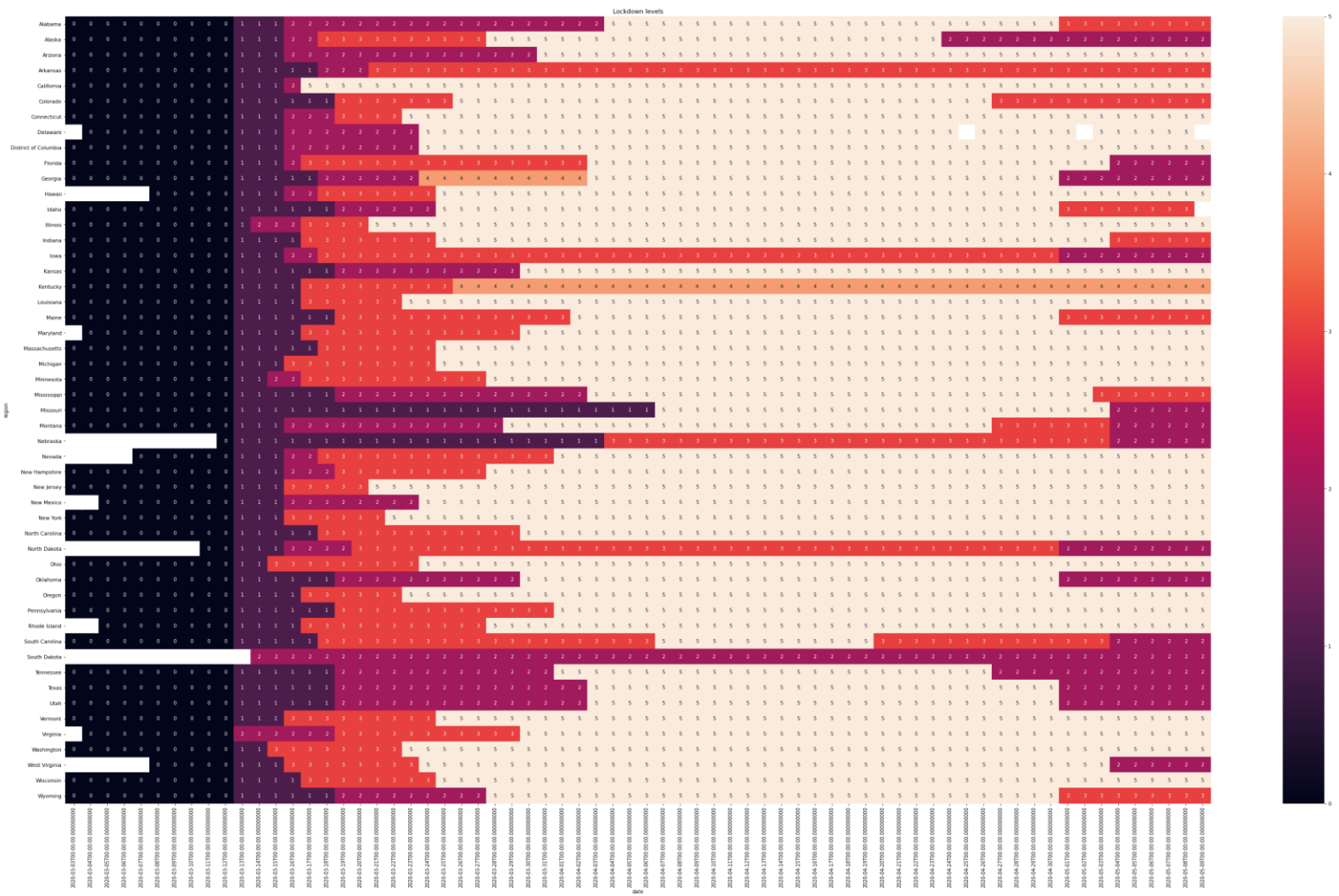
To produce the daily case numbers, I took both the 2nd order differences of the cumulative infected and their 7-day exponentially-weighted moving average. I took the differences of those, and used those as the infection acceleration target variable.

EDA:

I found that, for the vast majority of days, the R_t value was greater than 1.0:



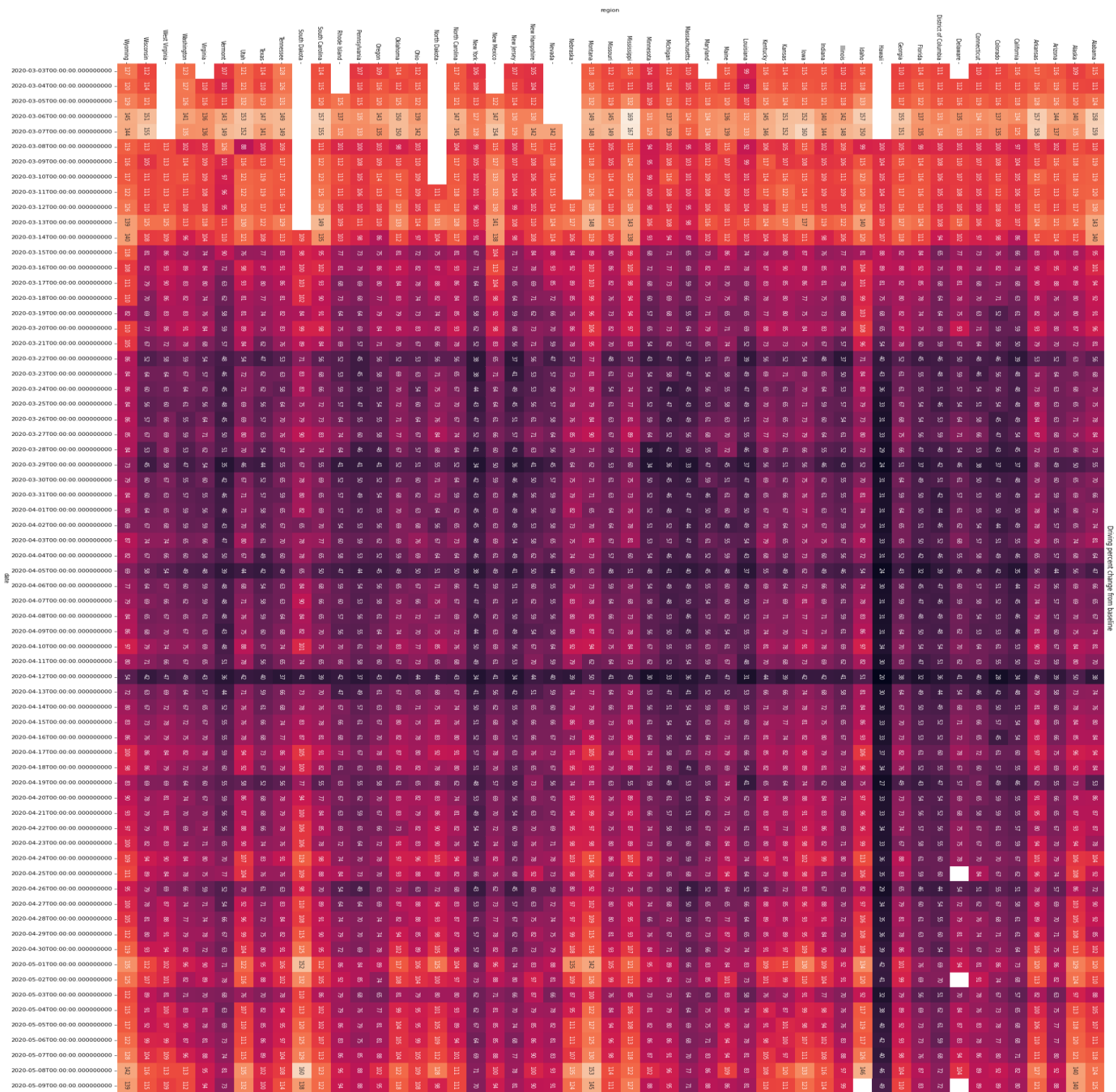
I was able to get 55 days of data spanning all states across all my datasets.



I plotted the lockdowns over time. The majority of states are at, or have been until very recently, a level 5 lockdown, though it is quite evident some states are reopening at the moment.

[illegible]

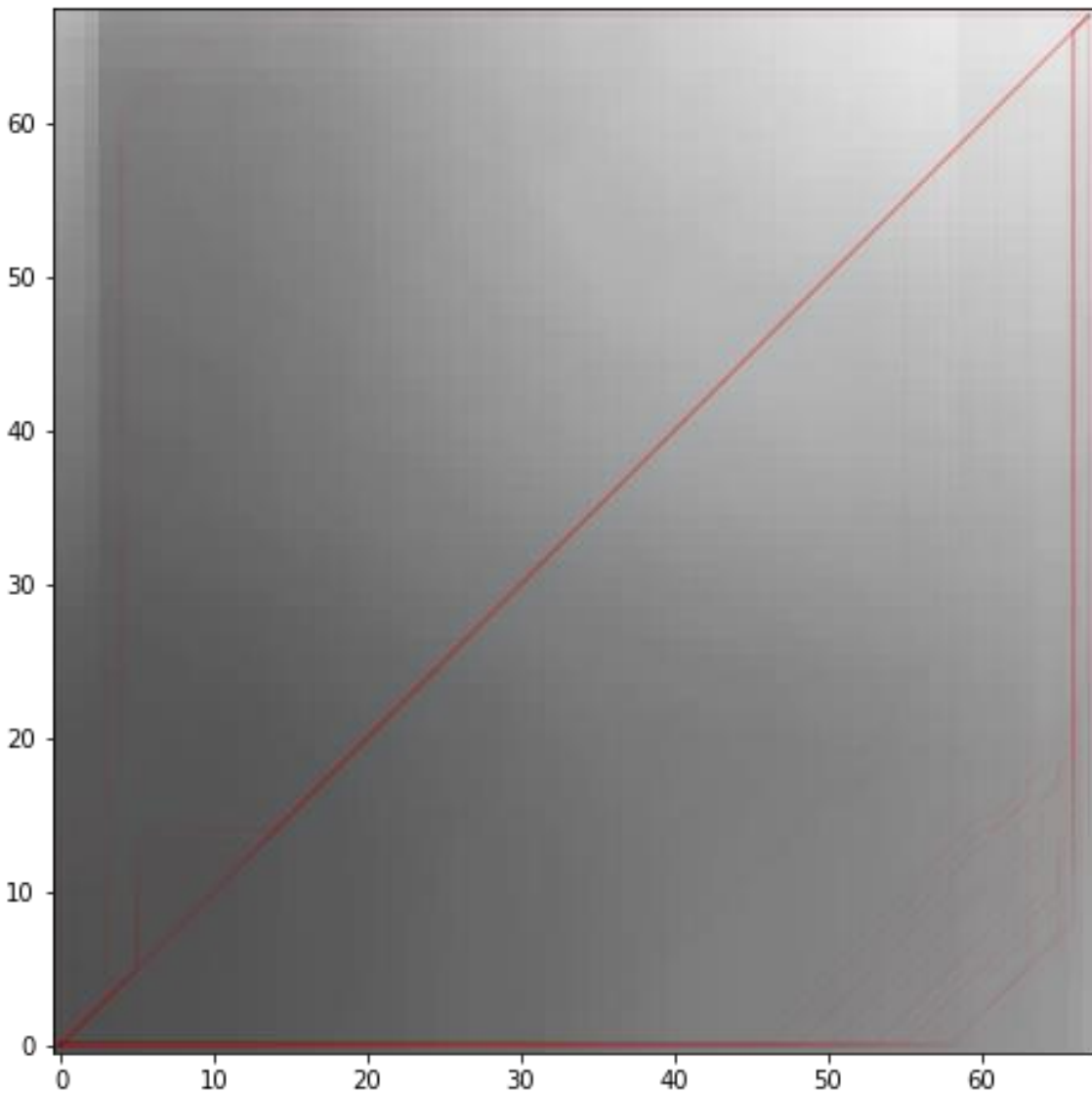
I also wanted to visualize the impact of the lockdowns on driving behavior.



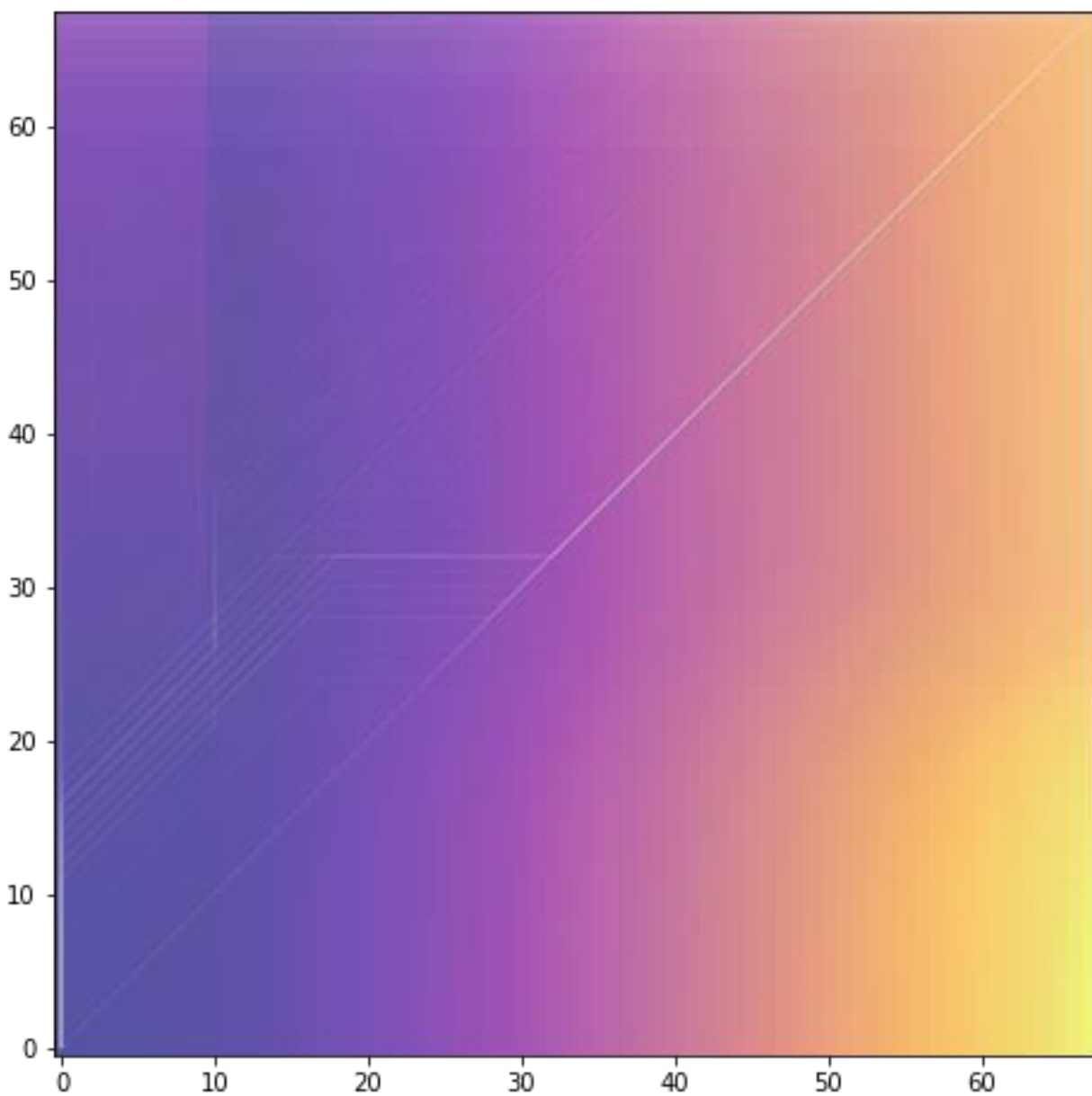
It is quite apparent from this that people significantly reduced their driving. It also seems people ramped up their driving in preparation of the lockdowns, driving quite significantly more than usual for a week or so before the lockdowns. It also seems that March 15th was the tipping point for driving percentage for the majority of states. It appears as though people were voluntarily altering their behavior before the harshest lockdown measures had arrived. Also, this heatmap clearly shows that April 12th was the lowest driving percentage day of all days of the crisis- April 12th, Easter Sunday this year, famously being President Trump's initially-hoped day for the re-opening of the country.

I also tried to do some dynamic-time-warping alignments between driving percent and daily cases, lockdown level and daily cases and driving percent and lockdown level, but they were fairly opaque for concluding things apart from that in the majority of cases, the best alignment of sequences was linear:

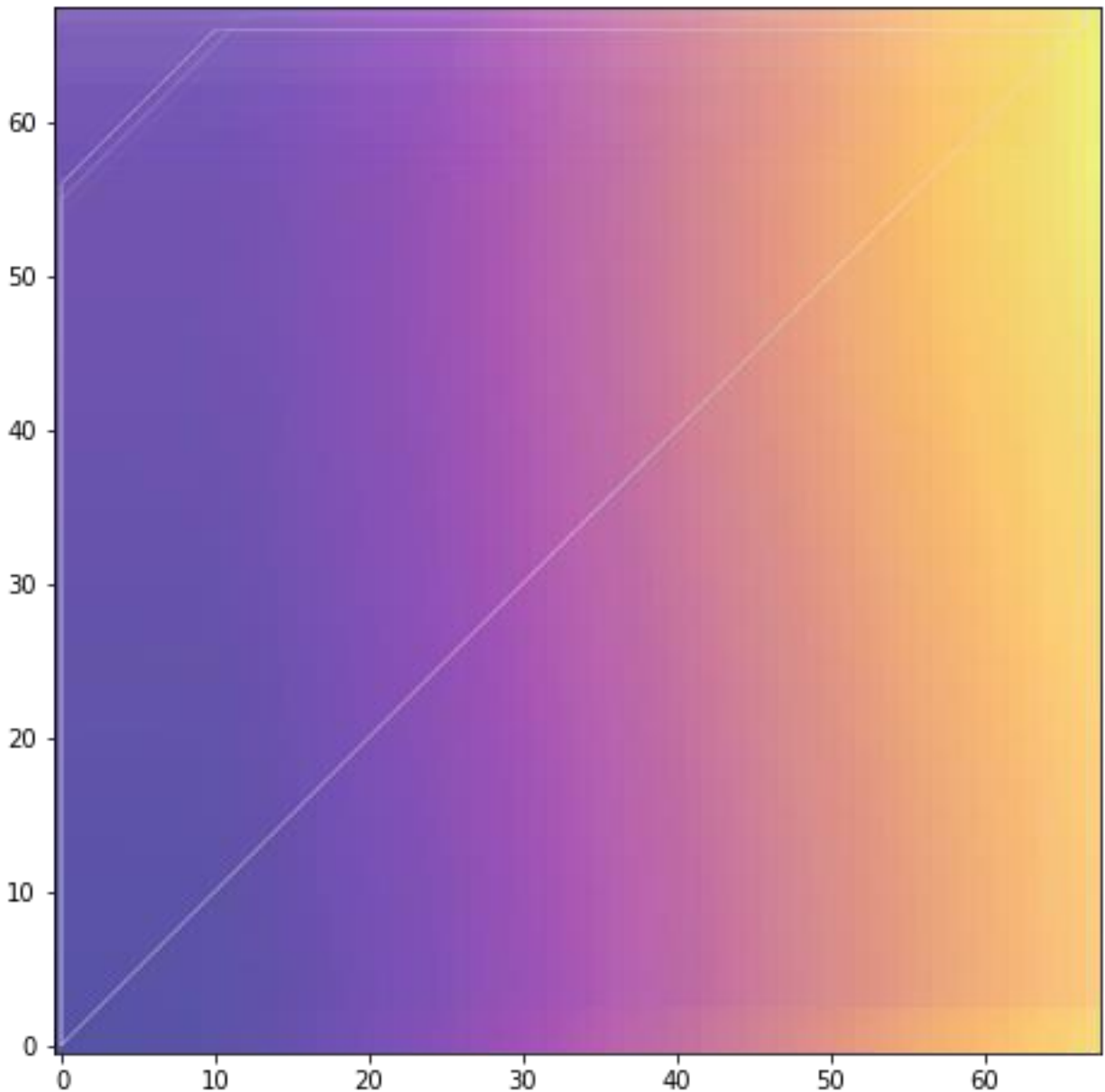
Driving percent (X) vs Daily Cases (Y) min-cost alignment curves (mean accumulated cost matrix is background)



Lockdown level vs daily cases alignment curves on top of mean accumulation matrix costs



Lockdown level vs driving percent mean alignment curves on top of mean accumulated cost of alignment



I wanted to see how good a model I could get from these data. At first I tried doing a regression for the R_t number, but this was a failure. I also tried to do an ordinal regression on R_t classes; e.g. predict 0 if $R_t < 0.8$, predict 1 if $0.8 < R_t < 1$, predict 2 if $1 < R_t < 1.2$, etc until predict 5 for $1.7 < R_t$, but this was not easy either.

Then I tried to do an LSTM model, which worked alright, but it wasn't that good. I did models at both the regional and the national level.

The last model that I came up with was a combination of a Temporal Convolutional Network (which I learned is a good architecture for timeseries) given the timeseries data as inputs, and another set of

inputs, the features that are not time-varying in my dataset (e.g. the obesity level as measured in 2010, or the urbanization percent for the state)

The inputs were as follows (all states' data were joined into one dataset):

Time-varying inputs: [

```
'driving_percent',  
'retail_and_recreation_percent_change_from_baseline',  
'grocery_and_pharmacy_percent_change_from_baseline',  
'transit_stations_percent_change_from_baseline',  
'workplaces_percent_change_from_baseline',  
'residential_percent_change_from_baseline',  
'lockdown_level_0',  
'lockdown_level_1',  
'lockdown_level_2',  
'lockdown_level_3',  
'lockdown_level_4',  
'lockdown_level_5',  
'lockdown_level_6',  
'lockdown_level_0_sum',  
'lockdown_level_1_sum',  
'lockdown_level_2_sum',  
'lockdown_level_3_sum',  
'lockdown_level_4_sum',  
'lockdown_level_5_sum',  
'lockdown_level_6_sum',  
'hospitalBedsRequired',  
'ICUBedsInUse',  
'positive_test_rate',  
'cumulativeDeaths', 'cumulativeInfected',  
'daily_cases',  
'daily_cases_ewm_avg',  
'RtIndicator',  
'cumulativePositiveTests', 'cumulativeNegativeTests'  
]
```

Non-time-varying inputs: ['lat', 'long',

```
'hospitalBedCapacity',  
'ICUBedCapacity',  
'poverty', 'age', 'income',  
'healthcare', 'healthcareLow', 'healthcareHigh',  
'obesity', 'smokes', 'smokesLow', 'smokesHigh',  
'Density per square mile of land area',  
'2018 Population',  
'urbanization']
```

The target variable was whether the daily cases was accelerating or not for that particular state.

The inputs 'lockdown_level_i' were indicator variables. I made lockdown level 0 equal to 1 if no lockdown measures were in place. If for e.g. lockdown level 4 were in place, I made lockdown level 0 equal 0, and made every lockdown_level_i for $1 \leq i \leq 4$ equal to 1. Lockdown_level_i_sum is the number of days that region has been at that lockdown level over time.

Here is my model architecture:

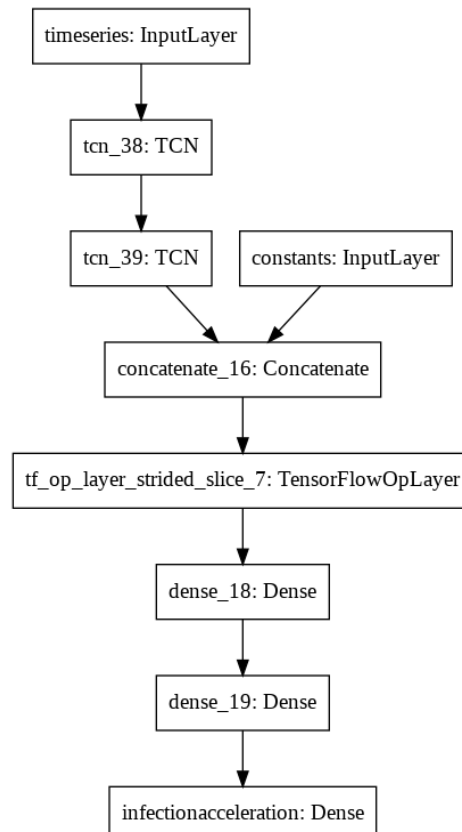
For temporal inputs:

1. TCN layer with kernel size = 3, dilations = [1,2,4,8,16,32], 2 residual stacks, 16 filters, returning sequences, fed into..
2. Another TCN layer with kernel size = 3, dilations = [1,2,4,8,16,32], 2 residual stacks, 16 filters, returning sequences

Then I concatenated the output of the last TCN layer (in sequence form) with the constant features (arranged as a constant sequence) into one tensor. Then with that vector I:

1. Took the last element of the sequence of concatenated vector (aka index[:, -1, :])
2. Fed it through 50-neuron fully-connected layer with ReLU activation
3. Fed it through 50-neuron fully-connected layer with ReLU activation again
4. Lastly ran through a 1-neuron fully-connected output layer with sigmoid activation

I used binary cross entropy to train the model.



I fed it sequences of 5 timesteps of data at a time. I held out the most recent 10 days' worth of data as a test set. At first, I used 10% of the train set for model validation purposes, but at test time then just used the entire training set. This meant roughly 45 days for training (~40 days training and 5 validation).

Of the 10 days of the test set, sequences of 5 days at a time were used for input, and 10-5-1 (total data minus sequence length minus 1 day because of predicting the output *after* the sequence input) = 4 days were predicted, which corresponded to May 4th through May 8th.

I used the Adam optimizer with a learning rate of 3e-3. I tracked the accuracy, precision, and recall per epoch, for 100 epochs. The precision is (true positive)/(actual results) and the recall is (true positive)/(predicted results).

It was important to me to have a high recall, because it would indicate the ratio of correct infection acceleration predictions to false infection deceleration predictions would be high. It is better to tell people it is not quite safe yet when it is actually safe, than to tell them it is now safe when it is actually dangerous.

Here are the performance metrics:

Lastly I grouped the states by their latest lockdown level. Per lockdown level, I determined which states were predicted to have decreasing or increasing daily cases.

Predictions:

```
{0: {0: [], 1: []},
 1: {0: [], 1: []},
 2: {0: ['South Carolina',
        'Oklahoma',
        'South Dakota',
        'Montana',
        'West Virginia',
        'Florida',
        'Georgia',
        'Alaska'],
     1: ['Texas',
        'North Dakota',
        'Missouri',
        'Nebraska',
        'Utah',
        'Iowa',
        'Tennessee']},
 3: {0: ['Arkansas', 'Alabama', 'Maine', 'Wyoming', 'Idaho', 'Mississippi'],
     1: ['Colorado', 'Indiana']},
 4: {0: [], 1: ['Kentucky']},
 5: {0: ['Oregon',
        'Hawaii',
        'New Jersey',
        'Vermont',
        'Nevada',
        'Rhode Island',
        'Massachusetts',
        'Michigan',
        'Connecticut'],
     1: ['Virginia',
```

```
'Illinois',
'Louisiana',
'Maryland',
'Pennsylvania',
'Kansas',
'Ohio',
'District of Columbia',
'Arizona',
'New York',
'Washington',
'Wisconsin',
'North Carolina',
'Delaware',
'Minnesota',
'California',
'New Mexico',
'New Hampshire']},
6: {0: [], 1: []}}
mean accuracy on test set: 0.7745098039215687
mean precision on test set: 0.39215686274509803
mean recall on test set: 0.39215686274509803
```

Unhappy about the recall, but I don't think this model was actually that bad.

According to <http://EndCoronavirus.org/states>, the states currently 'beating' COVID-19 are Alaska, Hawaii, Montana, Vermont, West Virginia and Wyoming. All of those states are listed as 0, or decreasing daily cases in my model. States that are listed as 'nearly there' are Idaho, Kansas, Maine, Michigan, New Jersey, and New York: my model puts Idaho, Maine, Michigan, and New Jersey as decreasing, but sees Kansas and New York as potentially increasing. I believe that due to the binary nature of 'decreasing or not' daily cases, I didn't account for states being far from their peak.