

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
 - 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

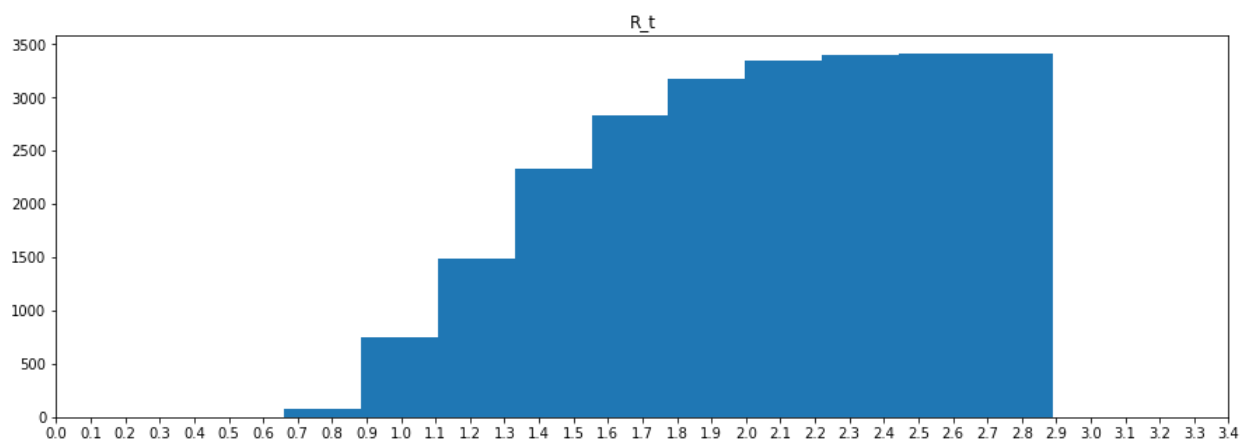
I also collected various non-time-series data per region:

- 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

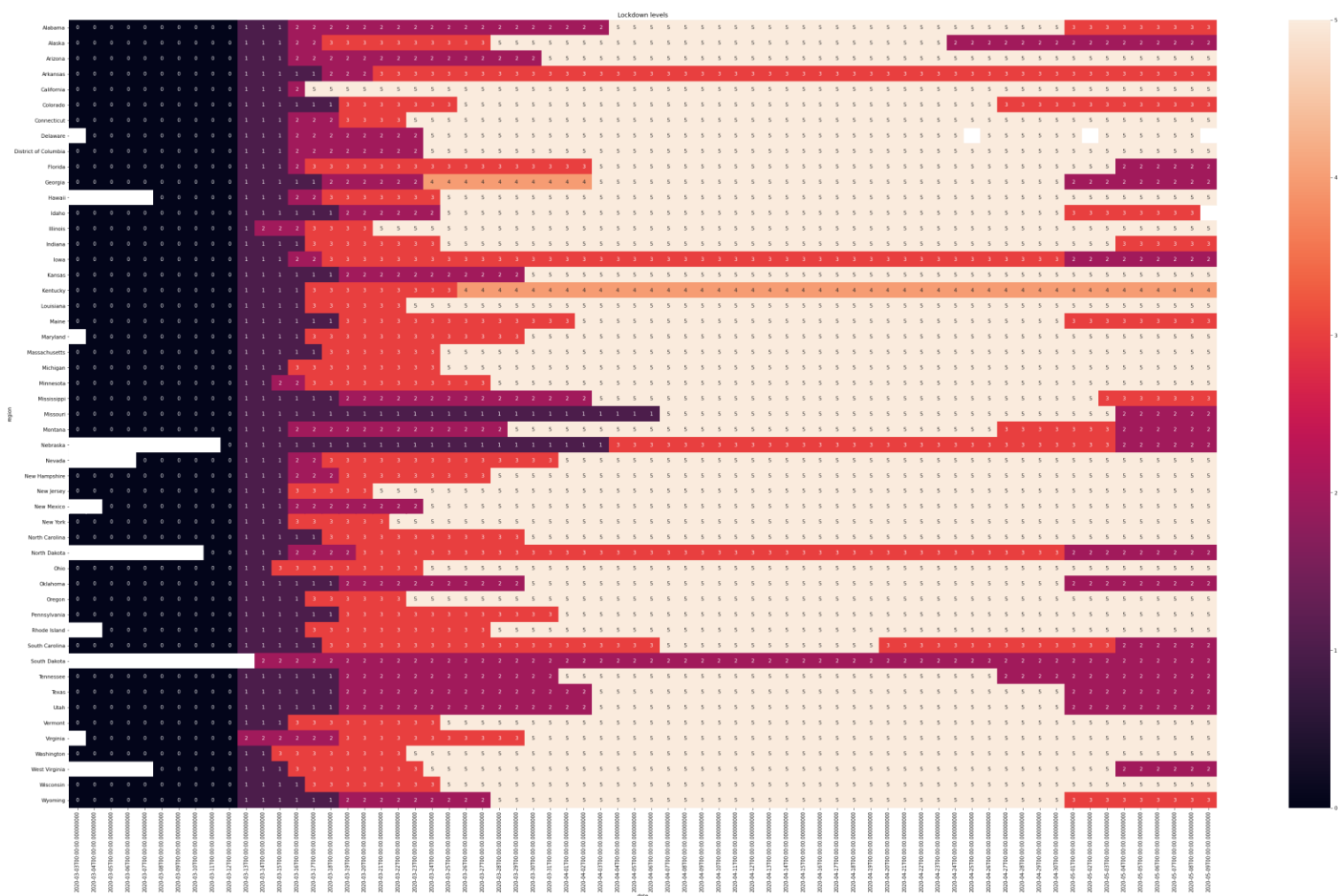
To produce the daily case numbers, I took both the 2nd order differences of the cumulative infected and their 7-day exponentially-weighted moving average. I took the differences of those, and used those as the infection acceleration target variable.

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.

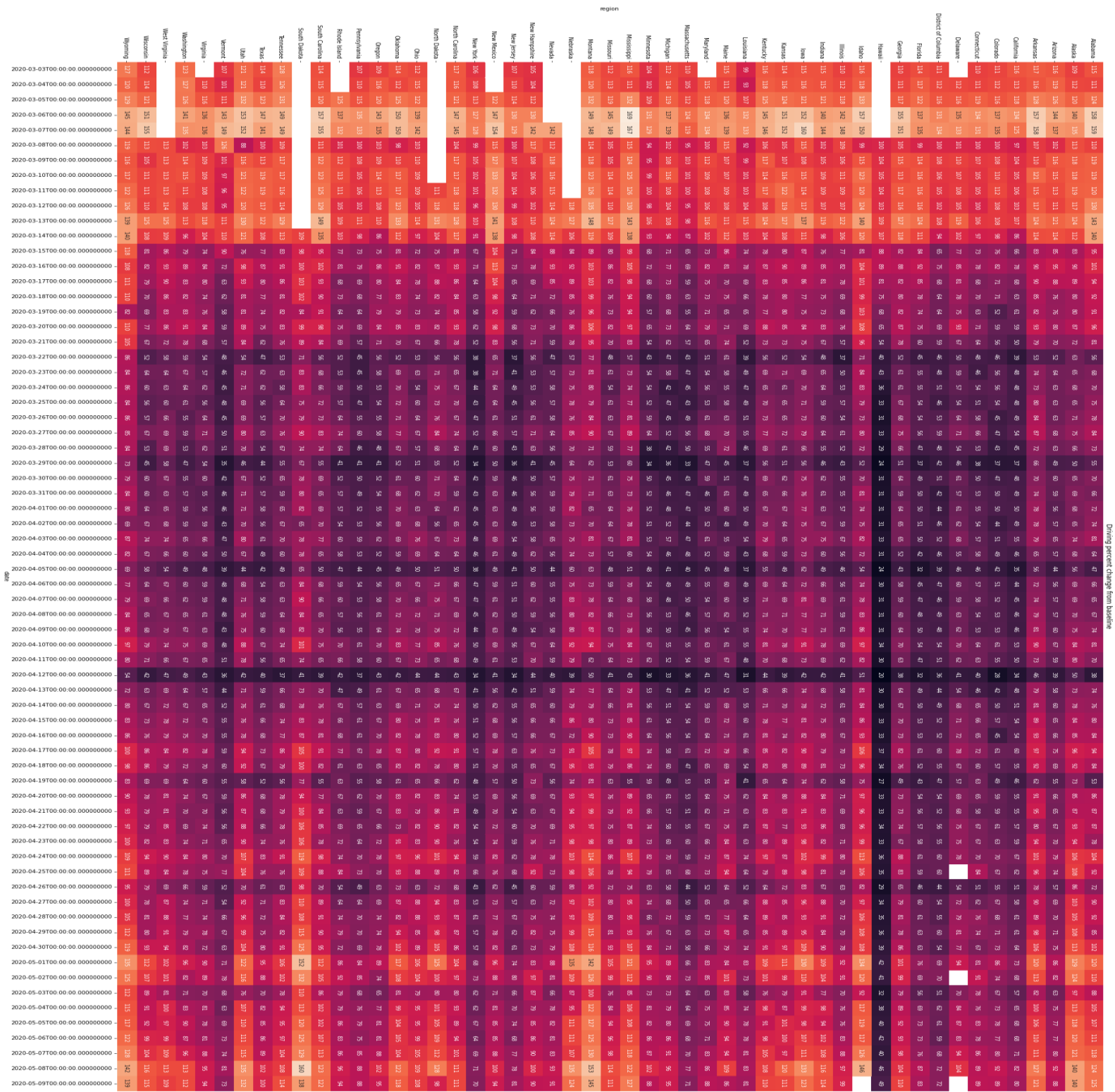


I plotted the lockdowns over time. The majority of states are at, or have been until very recently, a level 5 lockdown, though it is quite evident some states are reopening at the moment.

		Daily cases (normalized, 1=mean)
2020-04-29T00:00:00.0000000	Montana	0.16
2020-04-28T00:00:00.0000000	New York	0.24
2020-04-27T00:00:00.0000000	Alaska	0.28
2020-04-26T00:00:00.0000000	Michigan	0.32
2020-04-25T00:00:00.0000000	Vermont	0.34
2020-04-24T00:00:00.0000000	New Jersey	0.41
2020-04-23T00:00:00.0000000	Washington	0.46
2020-04-22T00:00:00.0000000	West Virginia	0.50
2020-04-21T00:00:00.0000000	Nevada	0.51
2020-04-20T00:00:00.0000000	Louisiana	0.55
2020-04-19T00:00:00.0000000	Wyoming	0.55
2020-04-18T00:00:00.0000000	Connecticut	0.57
2020-04-17T00:00:00.0000000	Florida	0.62
2020-04-16T00:00:00.0000000	Hawaii	0.63
2020-04-15T00:00:00.0000000	Pennsylvania	0.64
2020-04-14T00:00:00.0000000	Oklahoma	0.66
2020-04-13T00:00:00.0000000	Massachusetts	0.71
2020-04-12T00:00:00.0000000	Georgia	0.81
2020-04-11T00:00:00.0000000	South Dakota	0.82
2020-04-10T00:00:00.0000000	Oregon	0.84
2020-04-09T00:00:00.0000000	Arkansas	0.85
2020-04-08T00:00:00.0000000	Missouri	0.86
2020-04-07T00:00:00.0000000	South Carolina	0.86
2020-04-06T00:00:00.0000000	Colorado	0.87
2020-04-05T00:00:00.0000000	New Hampshire	0.88
2020-04-04T00:00:00.0000000	Maine	0.89
2020-04-03T00:00:00.0000000	Arizona	0.90
2020-04-02T00:00:00.0000000	District of Columbia	0.91
2020-04-01T00:00:00.0000000	Alabama	0.92
2020-03-31T00:00:00.0000000	California	0.93
2020-03-30T00:00:00.0000000	Delaware	0.94
2020-03-29T00:00:00.0000000	I Idaho	0.95
2020-03-28T00:00:00.0000000	Rhode Island	0.96
2020-03-27T00:00:00.0000000	Illinois	0.97
2020-03-26T00:00:00.0000000	Indiana	0.98
2020-03-25T00:00:00.0000000	Iowa	0.99
2020-03-24T00:00:00.0000000	Kansas	1.00
2020-03-23T00:00:00.0000000	Kentucky	1.01
2020-03-22T00:00:00.0000000	Maryland	1.02
2020-03-21T00:00:00.0000000	Minnesota	1.03
2020-03-20T00:00:00.0000000	Mississippi	1.04
2020-03-19T00:00:00.0000000	Nebraska	1.05
2020-03-18T00:00:00.0000000	New Mexico	1.06
2020-03-17T00:00:00.0000000	North Carolina	1.07
2020-03-16T00:00:00.0000000	North Dakota	1.08
2020-03-15T00:00:00.0000000	Ohio	1.09
2020-03-14T00:00:00.0000000	Tennessee	1.10
2020-03-13T00:00:00.0000000	Texas	1.11
2020-03-12T00:00:00.0000000	Utah	1.12
2020-03-11T00:00:00.0000000	Virginia	1.13
2020-03-10T00:00:00.0000000	Wisconsin	1.14

I plotted a heatmap of the daily cases over time, where each state's cases have been normalized to their maximum. As you can see, the vast majority of states are still above half of their max daily cases. Currently, 20 states are actually right now at their peak in daily cases.

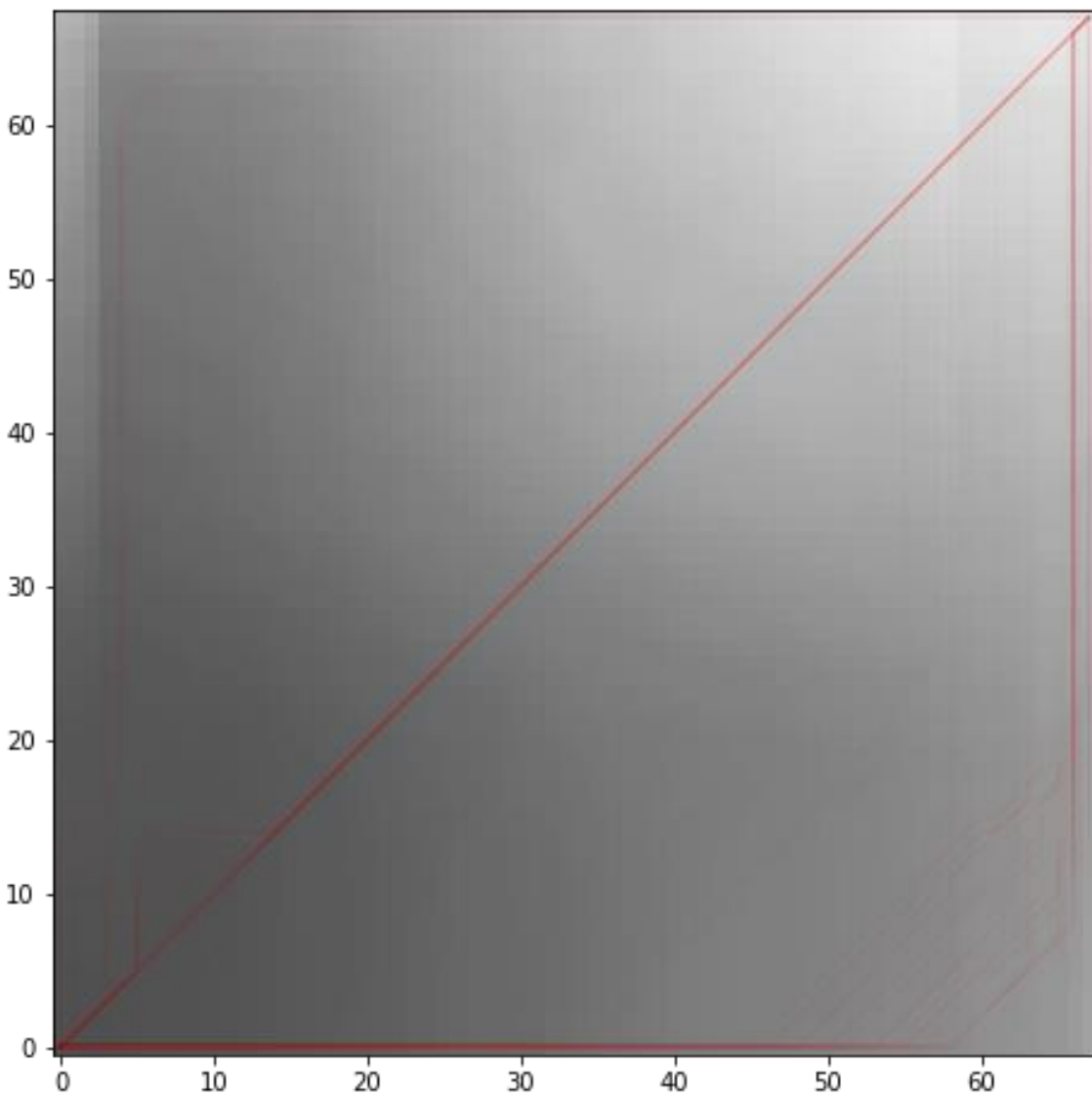
I also wanted to visualize the impact of the lockdowns on driving behavior.



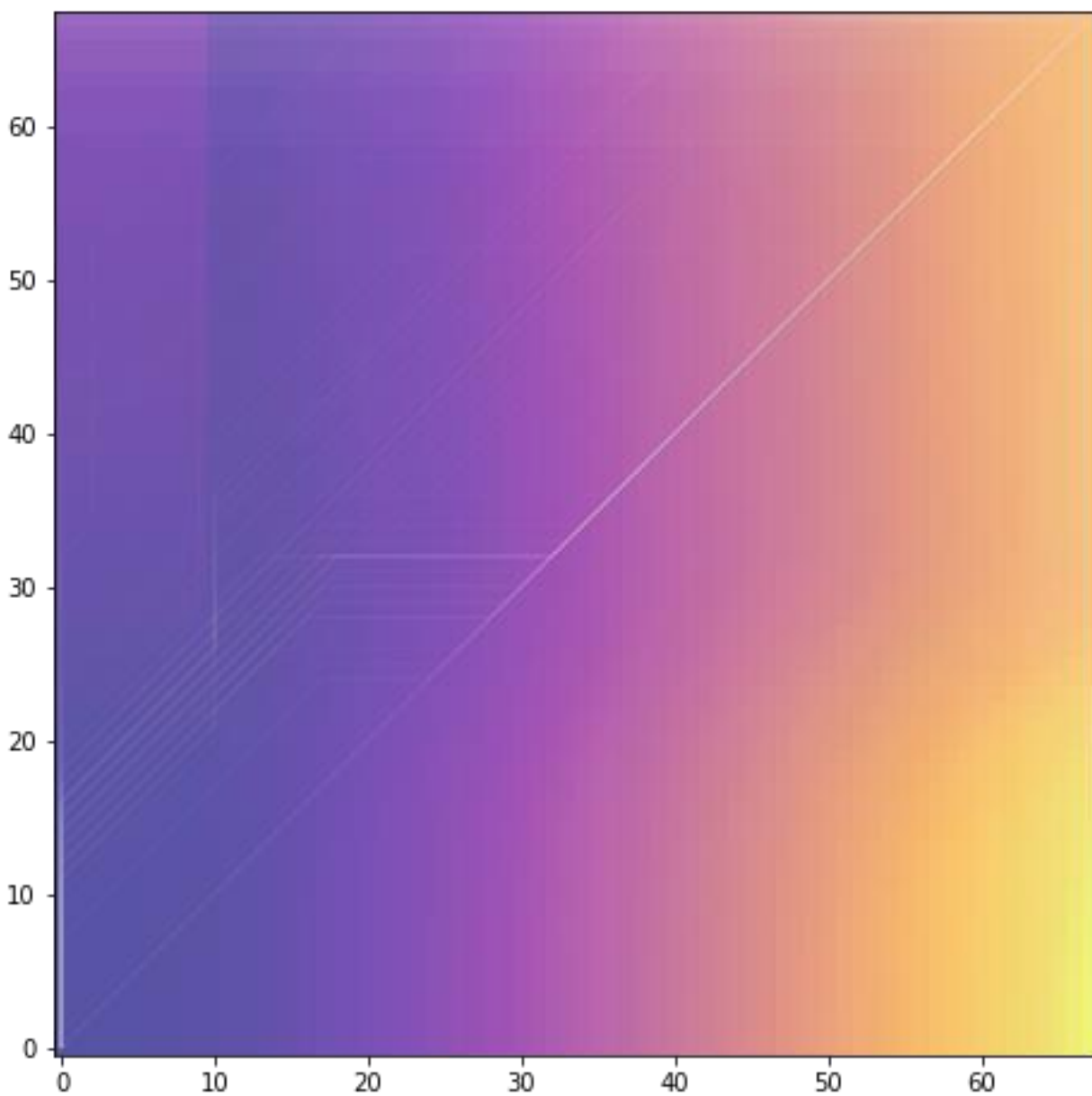
It is quite apparent from this that people significantly reduced their driving. It also seems people ramped up their driving in preparation of the lockdowns, driving quite significantly more than usual for a week or so before the lockdowns. It also seems that March 15th was the tipping point for driving 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:

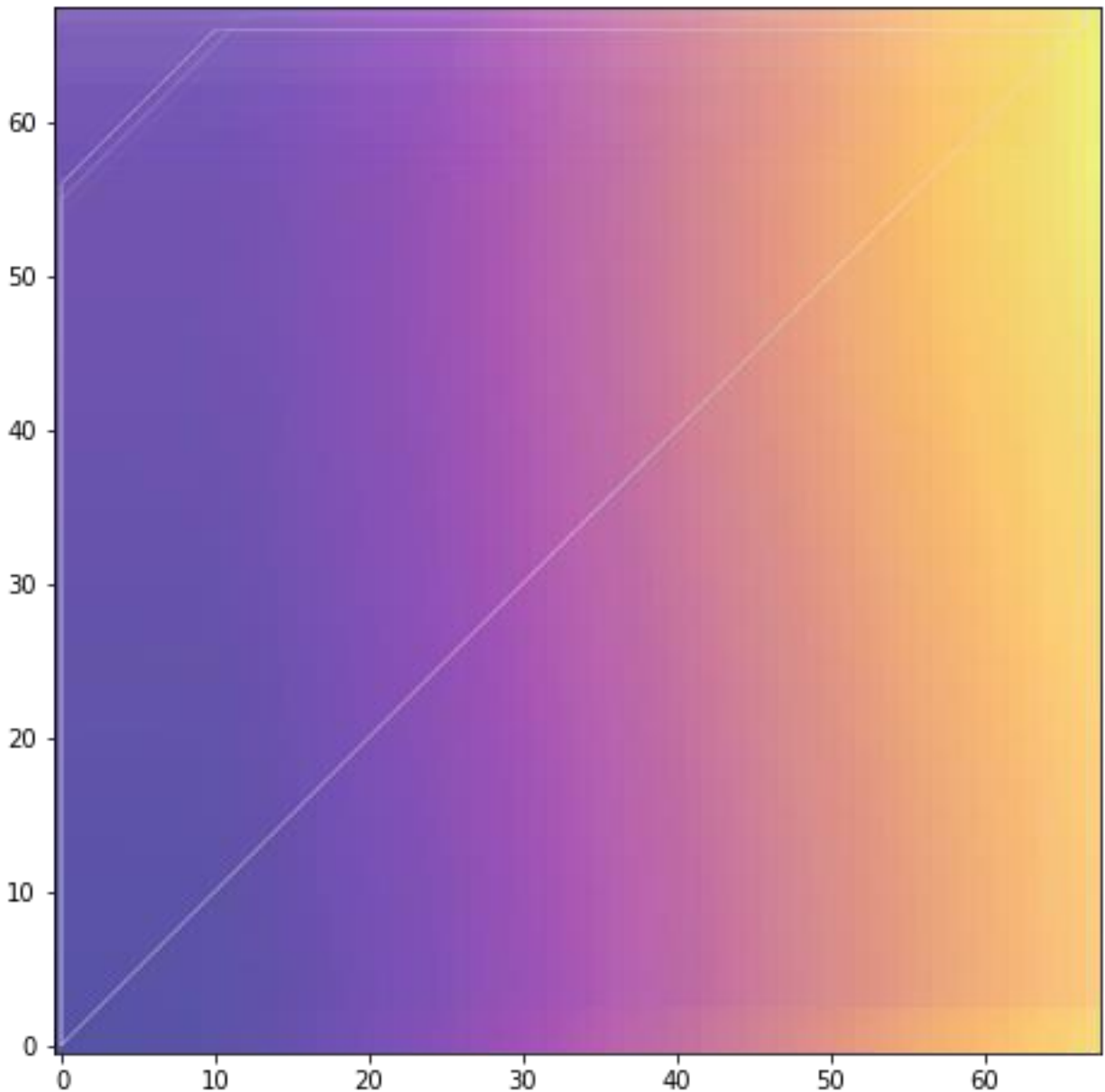
Driving percent (X) vs Daily Cases (Y) min-cost alignment curves (mean accumulated cost matrix is background)



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:

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',

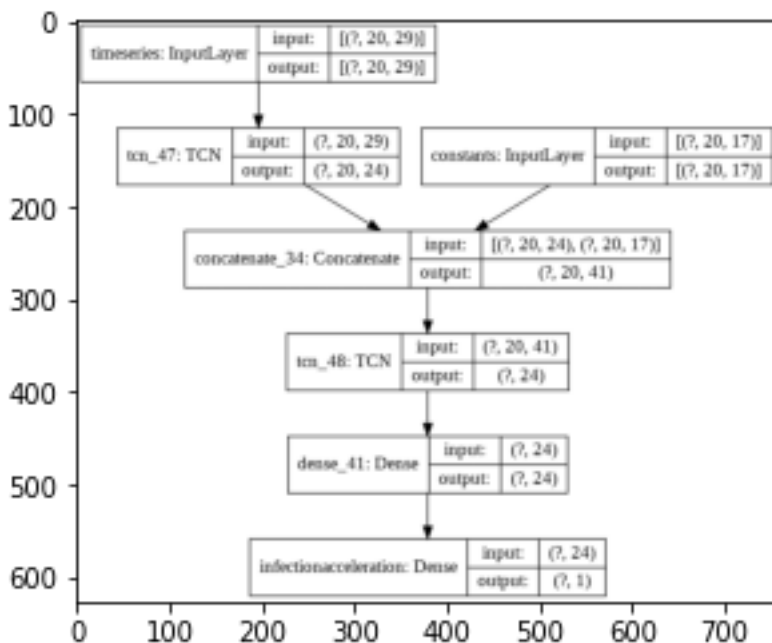
```
'hospitalBedCapacity',
'ICUBedCapacity',
'poverty', 'age', 'income',
'healthcare', 'healthcareLow', 'healthcareHigh',
'obesity', 'smokes', 'smokesLow', 'smokesHigh',
'Density per square mile of land area',
'2018 Population',
'urbanization']
```


The target variable was whether the last 7-days-exponentially-weighted daily cases was accelerating or not.

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.

I used 24 filters, kernel of size 2, and dilations of [1,2,4,8,16,32,64] for the first TCN layer. I made it return a sequence output, which I concatenated with the constant features. Then I fed both combined into another TCN layer that returned only a single output (ie, not a sequence) – I did this because I was having difficulty retrieving just the last element of the sequence returned from the TCN layer that returned sequences.. Then I fed this latest TCN layer output through a 24-node dense layer with a ReLU activation, which I then passed through another 1-node dense layer that had a sigmoid activation function.

I used binary cross entropy to train the model.



I fed it sequences of 20 timesteps of data at a time. I held out the most recent 25 days' worth of data as a test set. I used 10% of the train set to be a validation set. This meant roughly 27 days for training, and 3 for validation. Then of the 25 of the test set, 20 were used for data, and 25-20-1 (minus 1 day because of predicting the output *after* the sequence input) = 4 days were predicted.

I used the Adam optimizer with a learning rate of $1e-3$. 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.

By the 100th epoch, I had a training set accuracy of 93%, training set precision of 94% and training set precision of 97%. I had a validation set accuracy of 81%, precision 85%, and recall 90%. On the test set, given the 25 most recent days of data (previously unseen) in sequences of length 20, predicting the most recent 4 days was (averaged across all states) 53.4% accurate for all 4 days, 54% precise, and had a recall of 73% .

Then I grouped the states by their latest lockdown level. Of those in each level, I determined which were predicted to have a decelerating outcome, and which were predicted to have an accelerating outcome. There was an issue with my code (unless the outbreak is truly getting worse!), because for all regions in the USA my model predicted an acceleration of infections, for all lockdown levels.