

COVID-19 Forecasting And Analyzation

Weibang He, Yilin Zhou, Shuqing Zhang, Yuan Wan
Computer Engineering Department
San Jose State University, San Jose, California, USA

Abstract—In response to the COVID-19 pandemic, people from all kinds of science domain and industry domain work together and try their best to figure out the problem faced by all over the world. Global research has been taken to apply recent advances in machine learning techniques to generate new insights in support of the ongoing fight against this infectious disease. Even as young software engineering graduates in SJSU, we are trying to do our best in this challenging time. Here are what we do:

- 1) Overview including requirements, architecture, toolkit.
- 2) Select suitable machine learning algorithm on IBM Cloud and establish web UI showing the past COVID-19 situation and predicting the future positive cases based on each state in US.
- 3) Study the reason of COVID-19 outbreak in US and compare several possibilities.

I. OVERVIEW

A. Requirements

We want to build a ML model to predict the trend of the disease, which helps research on the future study of infectious diseases. And also, if accurate, this model can help people to prepare for COVID-19, such as how long they will stay at home in different areas, how many people can be infectious in the future, and which can be dangerous and which are not.

B. Architecture

-Data Set: We collect data from kaggle and github. The data source link is at the end of this file.

-Front End: React.js, Bootstrap. Diagram will be simpler and more intuitive to show data relation. As for chart showing, Echart.js is chosen to show our data.

-Back End: Python, Node.js, IBM cloud. We will collect plenty of data associated with COVID-19, then choose most suitable model, train our data

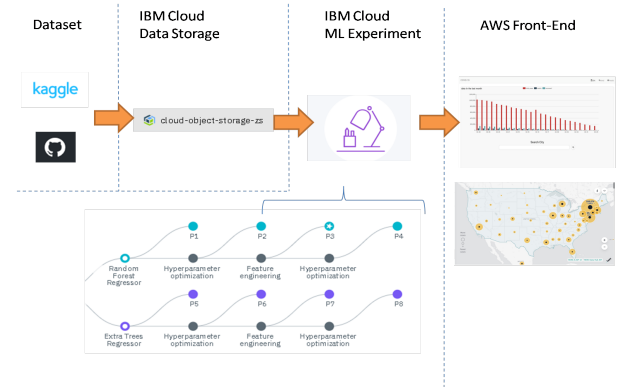


Fig. 1: architecture

in IBM Watson machine learning experiment, and provide an API for prediction.

We fetch the latest daily data about covid-19 in every state and the nearest one month in USA from Internet, then update them in our front page.

C. Tech Tools

- IBM Cloud Watson Machine Learning
- IBM Cloud Data Storage
- AWS
- ReactJS
- NodeJS
- Python

II. MODELING AND UI

A. Select suitable Machine Learning Algorithm on IBM Cloud and upload dataset for computing

We select (datetime, state) as the feature vector in the dataset and new positive population per day as prediction column. IBM Watson Studio provides auto AI development platform. What we need to do is to upload our dataset of COVID-19 to the platform. Then the platform runs and promotes several possible algorithms and calculate the RMSE.

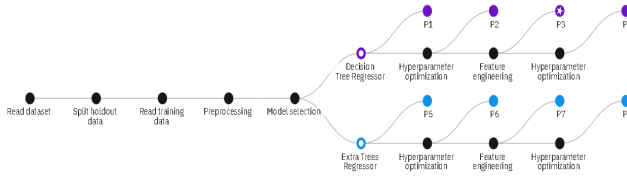


Fig. 2: SelectModel

The algorithm Decision Tree Regressor is selected, for which has the lowest RMSE of 1805. Then we generate a model from it. The Model Evaluation Measures are shown in Table I:

TABLE I: Model Evaluation Measures

	Holdout Score	Cross Validation Score
RMSE	1,292.887	1,805.030
R2	0.997	0.992
MSE	1,671,557.863	3,592,597.296

B. Establish Web UI showing the past COVID-19 situation and predicting the future positive cases based on each state in US

We establish a web site to offer the prediction of each state in US. Users can select the area which they are interested in and select the date they care about. The example is shown in Fig.3.

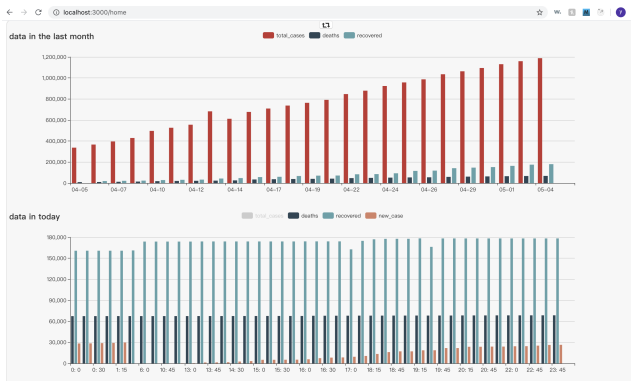


Fig. 3: predict

Also, we provide the current situation of COVID19, as map shown in Fig.4.

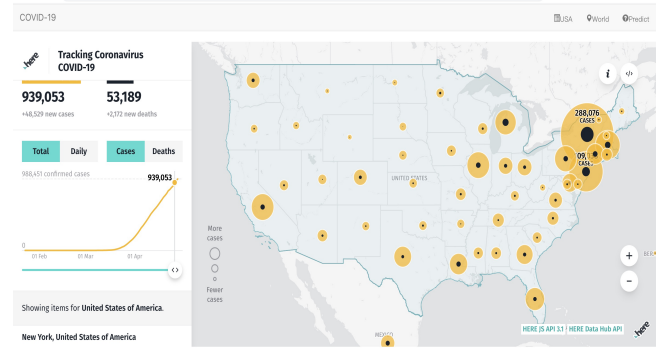


Fig. 4: map

III. STUDY THE REASON OF COVID-19 OUTBREAK IN US AND COMPARE SEVERAL POSSIBILITIES

A. Wearing Mask Influence on COVID-19

Is wearing a mask helpful for us to prevent from getting the COVID-19?

Does quarantine work for control/reduce the COVID-19 positive increase number for a country?

We all knew that China was the only country which implemented quarantine rule and compulsorily required every citizens to wear mask at the virus exploding period. China had most Covid-19 cases in all over the world, however now China already controlled and weaken Covid-19's spreading speed in China area and set up plan to resume work and resume cities' daily operation.

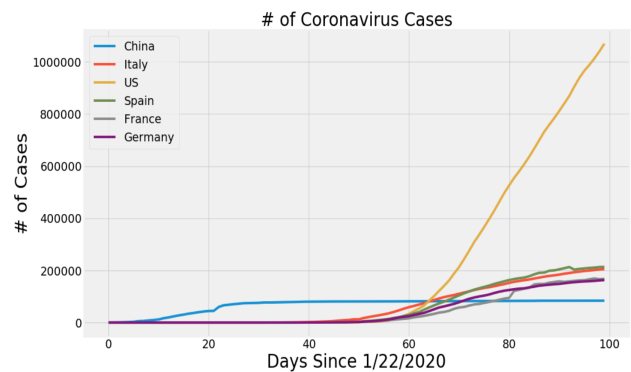


Fig. 5: nations

So, here are two questions: Is wearing mask helpful for us to prevent from getting the COVID-19 Virus? Does the quarantine rule works for control or weaken the Covid-19 positive cases' increase for a country?

We tracked back some news about the quarantine rule and government recommended citizens wear mask as necessary, just as Fig.6.

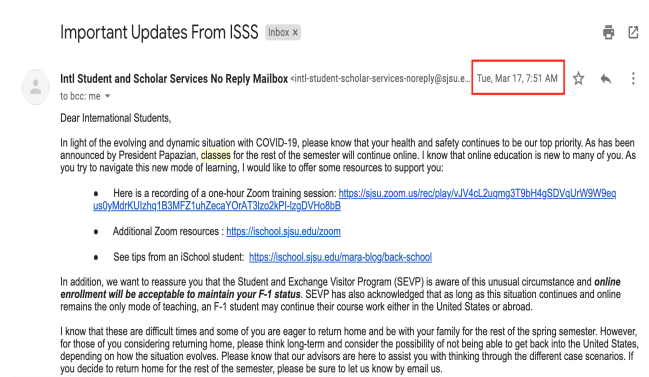


Fig. 6: ImportantUpdates

According to news' date, we could assume that the start date to implement the quarantine rule is Mar. 17th, and the date for CDC and government recommended people wearing facial masks is April 3rd. The dates around these two date (pre or post 2-3 days) also can be considered as the beginning of the implementation of rule suggestion. Then we will check whether data give us clues from these two dates.

Fig.7 and Fig.8 show the data about grocery & retailer mobility corresponding to date. The higher mobility explained the more people go to grocery & retailer compared to last day.

From the data, we can find that, the mobility increased dramatically from Mar.11th to Mar.17th, which also match the panic period in history—people purchased over amounts of grocery to stay at home for long time.

After Mar.17, we can find that the mobility begin to decrease, and the citizens did efficiently cooperate with government's rule.

To make the data more readable, we will fix the fluctuation on the data set. We will calculate the positive increase ratio per day. For example, Day 1, the increase number is 100, Day 2 is 50. Then the increase ratio should be -50%, showed in Fig.9.

According to the data set, we assume that the COVID-19 explode date for the US nation should be on Feb 28th. Before Feb 28th, the number of COVID-19 case still a few. Then we can skip till Feb 28th. This is showed in Fig.10

	location	loc_type	parent_loc	mobility_type	date	mobility_change
0	Afghanistan	parent	world	Retail & recreation	2020-02-16 00:00:00	0.028430
1	Afghanistan	parent	world	Retail & recreation	2020-02-17 00:00:00	0.057307
2	Afghanistan	parent	world	Retail & recreation	2020-02-18 00:00:00	0.025352
3	Afghanistan	parent	world	Retail & recreation	2020-02-19 00:00:00	-0.008219
4	Afghanistan	parent	world	Retail & recreation	2020-02-20 00:00:00	-0.017038
...
942434	Weston County	child	Wyoming	Residential	2020-04-06 00:00:00	-0.387755
942435	Weston County	child	Wyoming	Residential	2020-04-07 00:00:00	-0.367347
942436	Weston County	child	Wyoming	Residential	2020-04-08 00:00:00	-0.360000
942437	Weston County	child	Wyoming	Residential	2020-04-09 00:00:00	-0.312500
942438	Weston County	child	Wyoming	Residential	2020-04-10 00:00:00	-0.431818

Fig. 7: retail0

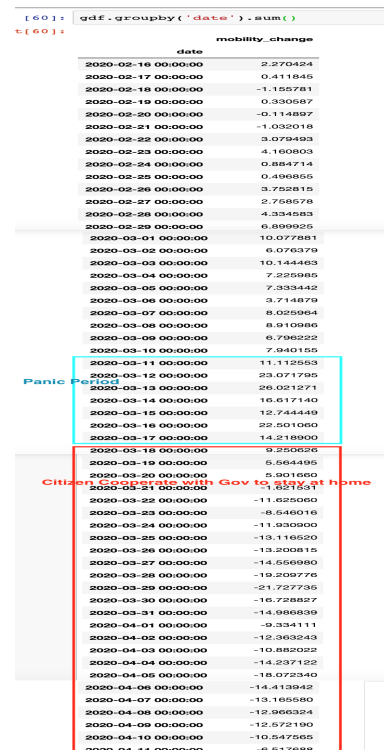


Fig. 8: retail1

The trend of increase rate is not stable, it will be influenced by each day's increase cases dramatically. It would be better for us to count the average to get a period's increase rate. Then we can compare the data between the period before quarantine and after the quarantine, and the period before promoting mask versus after promoting mask.

And we can view that, after the government implements quarantine rule, the increase positive ratio got obvious decrease. From around 45% down to 15%. We can view that the COVID-19

```
def2['recovered','positiveIncrease']
def3['increase_rate']=p[0:1]
for q in range(1,len(def3)):
    def3['increase_rate'][q-1]=(def3['positiveIncrease'][q]-def3['positiveIncrease'][q-1])/100
def3
```

	positive	positiveIncrease	recovered	death	day	increase_rate
0122	1.0	0.0	0.0	0.0	1	NaN
0123	1.0	0.0	0.0	0.0	2	NaN
0124	1.0	0.0	0.0	0.0	3	NaN
0125	1.0	0.0	0.0	0.0	4	NaN
0126	1.0	0.0	0.0	0.0	5	NaN
0127	1.0	0.0	0.0	0.0	6	NaN
0128	1.0	0.0	0.0	0.0	7	NaN
0129	1.0	0.0	0.0	0.0	8	NaN
0130	1.0	0.0	0.0	0.0	9	NaN
0131	1.0	0.0	0.0	0.0	10	NaN
0201	1.0	0.0	0.0	0.0	11	NaN
0202	1.0	0.0	0.0	0.0	12	NaN
0203	1.0	0.0	0.0	0.0	13	NaN
0204	1.0	0.0	0.0	0.0	14	NaN
0205	1.0	0.0	0.0	0.0	15	NaN
0206	1.0	0.0	0.0	0.0	16	NaN
0207	2.0	0.0	0.0	0.0	36	NaN
0208	2.0	0.0	0.0	0.0	37	NaN
0209	9.0	7.0	0.0	4.0	38	38.57429
0210	16.0	8.0	0.0	8.0	39	33.33333
0211	40.0	12.0	0.0	8.0	40	8.333333
0212	53.0	13.0	0.0	11.0	41	216.304615
0213	84.0	41.0	0.0	14.0	42	-12.180122
0214	207.0	36.0	0.0	16.0	43	85.555556
0215	273.0	66.0	0.0	20.0	44	67.862500
0216	387.0	100.0	0.0	26.0	45	35.778117
0217	538.0	148.0	0.0	27.0	46	23.648649
0218	721.0	183.0	0.0	31.0	47	26.05814
0219	1013.0	292.0	0.0	35.0	48	-6.548164
0220	1380.0	367.0	0.0	37.0	49	-46.816479
0221	1672.0	360.0	0.0	43.0	50	17.867143

Fig. 9: positiveIncrease

	positive	positiveIncrease	recovered	death	day	increase_rate	avg3_count
0222	9.0	7.0	0.0	4.0	38	28.57429	54.685518
0223	18.0	9.0	0.0	6.0	39	33.33333	54.685518
0224	40.0	12.0	0.0	8.0	40	8.333333	54.685518
0225	53.0	13.0	0.0	11.0	41	216.304615	54.685518
0226	84.0	41.0	0.0	14.0	42	-12.180122	54.685518
0227	207.0	36.0	0.0	16.0	43	85.555556	53.447834
0228	273.0	66.0	0.0	20.0	44	67.862500	53.447834
0229	387.0	100.0	0.0	26.0	45	35.778117	53.447834
0230	538.0	148.0	0.0	27.0	46	23.648649	53.447834
0231	721.0	183.0	0.0	31.0	47	26.05814	53.447834
0232	1013.0	292.0	0.0	35.0	48	-6.548164	26.05814
0233	1380.0	367.0	0.0	37.0	49	-46.816479	26.05814
0234	1672.0	360.0	0.0	43.0	50	17.867143	42.528621
0235	4943.0	1084.0	0.0	77.0	54	17.598979	42.528621
0236	6732.0	1738.0	0.0	86.0	55	68.367305	42.528621
0237	8388.0	2158.0	0.0	130.0	56	18.418126	42.528621
0238	10865.0	2577.0	0.0	145.0	57	62.302606	42.528621
0239	15903.0	4180.0	0.0	180.0	58	37.813149	30.447707
0240	24816.0	9177.0	0.0	203.0	59	13.388587	30.447707
0241	37372.0	9037.0	0.0	366.0	60	37.201470	30.447707
0242	36334.0	6862.0	0.0	436.0	61	18.158670	30.447707
0243	47913.0	10679.0	0.0	521.0	62	-4.802821	30.447707
0244	57179.0	10160.0	0.0	720.0	63	30.383201	14.783517
0245	69473.0	12094.0	147.0	863.0	64	40.848195	14.783517
0246	86789.0	17316.0	87.0	1231.0	65	7.836663	14.783517
0247	105482.0	16673.0	242.0	1604.0	66	5.641022	14.783517
0248	124815.0	19025.0	314.0	2205.0	67	0.590503	14.783517
0249	144207.0	20482.0	4081.0	2927.0	68	8.941587	10.53332
0250	160521.0	21224.0	4560.0	3900.0	69	15.328988	10.53332
0251	188882.0	24877.0	5666.0	3877.0	70	2.861986	10.53332
0252	215717.0	25179.0	7084.0	4803.0	71	11.434130	10.53332
0253	243335.0	26556.0	8568.0	5824.0	72	14.045505	10.53332
0254	270234.0	31890.0	10891.0	7112.0	73	4.747223	-0.410327
0255	306712.0	38116.0	13642.0	8472.0	74	16.341177	-0.410327
0256	334718.0	35966.0	14563.0	9714.0	75	10.710158	-0.410327
0257	363465.0	38747.0	16884.0	10837.0	76	5.781473	-0.410327
0258	392817.0	36409.0	18471.0	12811.0	77	-3.763262	-0.410327
0259	423645.0	38773.0	21163.0	14737.0	78	13.453589	-3.763262
0260	456260.0	34215.0	24888.0	16658.0	79	0.967412	-3.763262
0261	490836.0	34546.0	28054.0	18731.0	80	-13.502734	-3.763262
0262	522837.0	26891.0	31821.0	20880.0	81	-2.274455	-3.763262
0263	551906.0	29219.0	34913.0	22237.0	82	-14.543320	-3.763262
0264	578475.0	24969.0	35442.0	23764.0	83	3.300292	2.726484
0265	609688.0	29793.0	36537.0	26065.0	84	16.032024	2.726484
0266	633112.0	30444.0	40322.0	28524.0	85	1.882144	2.726484
0267	661428.0	31017.0	44840.0	30732.0	86	1.982794	2.726484
0268	690761.0	31832.0	49364.0	32785.0	87	-11.438622	2.726484
0269	720717.0	38115.0	60981.0	34873.0	88	-1.796356	2.693403
0270	751287.0	37511.0	67338.0	38224.0	89	-6.647453	2.693403
0271	779419.0	28132.0	69636.0	37913.0	90	3.098034	2.693403
0272	803260.0	29911.0	73002.0	40471.0	91	7.648261	2.693403
0273	830233.0	27893.0	76303.0	43058.0	92	16.164359	2.693403
0274	861788.0	31955.0	82194.0	44583.0	93	9.355296	-3.716290
0275	895006.0	34518.0	101517.0	46251.0	94	4.108002	-3.716290
0276	932242.0	38836.0	112783.0	48989.0	95	-24.821908	-3.716290
0277	960226.0	27046.0	118871.0	49716.0	96	-16.020763	-3.716290
0278	981134.0	21876.0	121608.0	50237.0	97	11.802889	-3.716290
0279	1005882.0	24458.0	138342.0	52025.0	98	12.702410	4.693114
0280	1031017.0	27965.0	147484.0	53225.0	99	1.374032	4.693114
0281	1061101.0	27845.0	168957.0	52965.0	100	0.000000	4.693114

Fig. 10: Feb28th

positive case increase speed gets slow down.

More importantly, after April 3rd, the COVID-19 positive case increase ration from 15% down to below 5%. It is very convincing that quarantine rule/promoting mask have positive effect for a country to fight with COVID-19.

B. US doesn't pass the turning point

The negative increase of the COVID-19 positive case does not mean the US nation already pass the

turning point. Actually, the US nation still have not passed the turning point, though some state already passed the turning point.

We have implemented a turning point algorithm. Let us just take the Italy as an example.

	CurrentPositiveCases	NewPositiveCases	day	avg3_count	avg5-count
02-24	221	221	1	130.666667	176.0
02-25	311	93	2	130.666667	176.0
02-26	385	78	3	130.666667	176.0
02-27	588	250	4	242.666667	176.0
02-28	821	238	5	242.666667	176.0
02-29	1049	240	6	242.666667	440.2
03-01	1577	566	7	458.000000	440.2
03-02	1835	342	8	458.000000	440.2
03-03	2263	466	9	458.000000	440.2
03-04	2706	587	10	711.333333	440.2
03-05	3296	769	11	711.333333	1216.6
03-06	3916	778	12	711.333333	1216.6
03-07	5061	1247	13	1512.000000	1216.6
03-08	6387	1492	14	1512.000000	1216.6
03-09	7985	1797	15	1512.000000	1216.6

Fig. 11: Italytp

Because the new increase case get influenced every day, we can count the average increase cases for a period first. We can set the ratio as 3-day a period and 5-day a period. 3-day period could get more in precise date for turning point. However, if the data fluctuates in very huge range, we would recommend to use large scale range.

```
if ee['avg5-count'][m]==ee['avg5-count'][m-1]:
    Dlist.append(m-1)
    continue
else:
    Dlist.append(m-1)
    break
break
ThreeDlist = []
for x in range(1,len(Dlist)):
    if ee['avg3-count'][Dlist[x-1]]<ee['avg3-count'][Dlist[x]]:
        continue
    else:
        ThreeDlist.append(Dlist[x-1])
        ThreeDlist.append(Dlist[x-1]-1)
        ThreeDlist.append(Dlist[x-1]-2)
    for s in range(len(ThreeDlist)):
        print(ThreeDlist[s])
32
31
30
```

/usr/local/lib/python3.7/site-packages/ipykernel_launcher.py:6: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

Fig. 12: tp

From the turning point algorithm, we know that, the turning point will be appear around the 30-32 row of the dataframe(Fig.12), which is between the date from 3-25 to 3-27.

And we do find it the increase case vector get decrease after the turning point. (Fig.13)

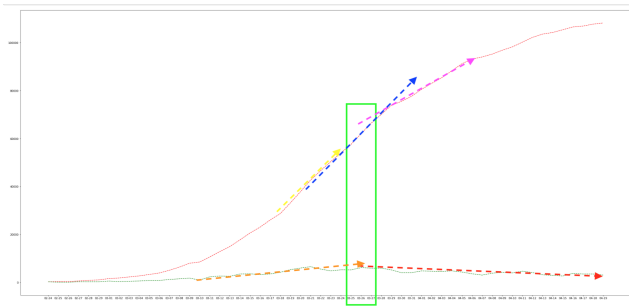


Fig. 13: Line chart

IF we implement this into the US dataset:(Fig.14)

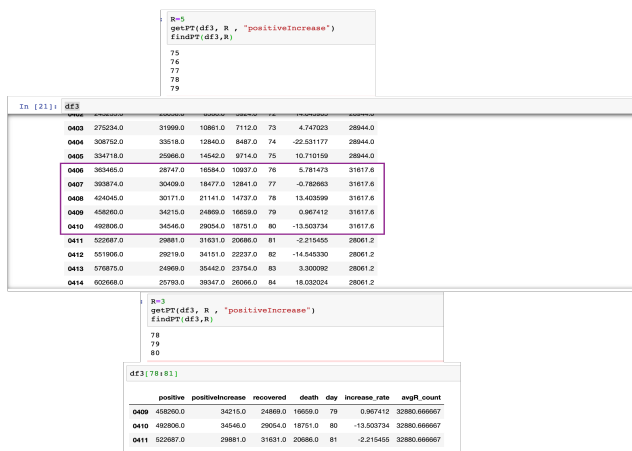


Fig. 14: US datasheet

We will find it the turning point for the US nation not come yet, but at the period from April 6 to April 10 will have the most increase case per day. After that, later new COVID-19 cases per day have negative increase with very low rate.

Why quarantine rule and wearing mask work, why we still have lot of increase case per day? Previous work, such as the quarantine rule and wearing mask, we just control the increase rate for new COVID-19 cases. If we do not implement those works, right now, we may have not rather 30000 new cases per day, but 50000, 60000 new cases per day.

Why we still don't have turning point yet? Although some states already passed the turning point, the U.S is a very big country, that COVID-19 condition is different from states to states. Until all the states from U.S pass the turning point, then the whole nation will pass the turning point then.

IV. CONCLUSIONS

In conclusion, we just control and reduce the accelerate rate for new cases per day, however there is still a distance for us to ultimately beat the virus for the whole nation.

REFERENCES

- [1] <https://www.politico.com/interactives/2020/coronavirus-testing-by-state-chart-of-new-cases/>
- [2] USA data set source: <https://www.kaggle.com/sudalairajkumar/covid19-in-usa>
- [3] <https://covid-19.direct/county/CA/Santa%20Clara>
- [4] <https://aqicn.org/data-platform/register/>
- [5] Fig.5 nations: <https://www.kaggle.com/tarunkr/covid-19-case-study-analysis-viz-comparisons>
- [6] Mobility data set source: <https://www.kaggle.com/lanheken0/community-mobility-data>
- [7] Italy Covid-19 data set source: https://www.kaggle.com/sudalairajkumar/covid-in-italy#covid19_italy_region.csv