

## 1. Data Cleaning

- - The given data was the cumulative cases and death in NV and NM.
- - We first obtained the daily data from this to better analyse the data cleaning and inferences.
- 
- - We also noticed that for some days the value of the covid cases/deaths were 0. This can be acceptable for the initial months when the impact of COVID was minimal. But after the first few months, the cases and deaths started rising. Getting a 0 value for a day, especially when the trend from the data clearly shows that value lies in a range of tens or hundreds (or even thousands), can be interpreted as missing values rather than the true figure.

These missing values were also replaced by previous days' values in such a manner that trend is not disturbed and data's consistency is maintained.

- - Next, we applied the Tukey's rule to find the outliers and remove them accordingly. We faced 2 major issues with applying Tukey's rule on the whole dataset:

1) As the COVID date is increasing (likely in a geometric distribution manner), and the cases are very less in the initial months and are very high in the last few months, if the Tukey's rule is applied on the whole dataset at once, then most of the points of the final few months will be classified as outliers.

2) Also, if there is a high value in the first few months or a low value in the final few months, then it should be flagged. But if Tukey's rule is applied on the whole data at once, this case won't be flagged.

To deal with this issue, we applied Tukey's rule on periods of 30 days each. This prevented the values in final few months to be marked as outliers. It also helped in removing unusually high or low values from each 30 day period.

Here are the outliers in the data:

- - The cumulative cases and deaths were calculated from this cleaned data and added back in.

### Normal Data -

	Date	NM confirmed	NV confirmed	NM deaths	NV deaths
0	2020-01-22	0	0	0	0
1	2020-01-23	0	0	0	0
2	2020-01-24	0	0	0	0
3	2020-01-25	0	0	0	0
4	2020-01-26	0	0	0	0
..	...	...	...	...	...
433	2021-03-30	191380	298052	3932	5127
434	2021-03-31	191655	298328	3937	5137
435	2021-04-01	191948	298651	3942	5144
436	2021-04-02	192156	298651	3949	5144
437	2021-04-03	192156	299440	3949	5161

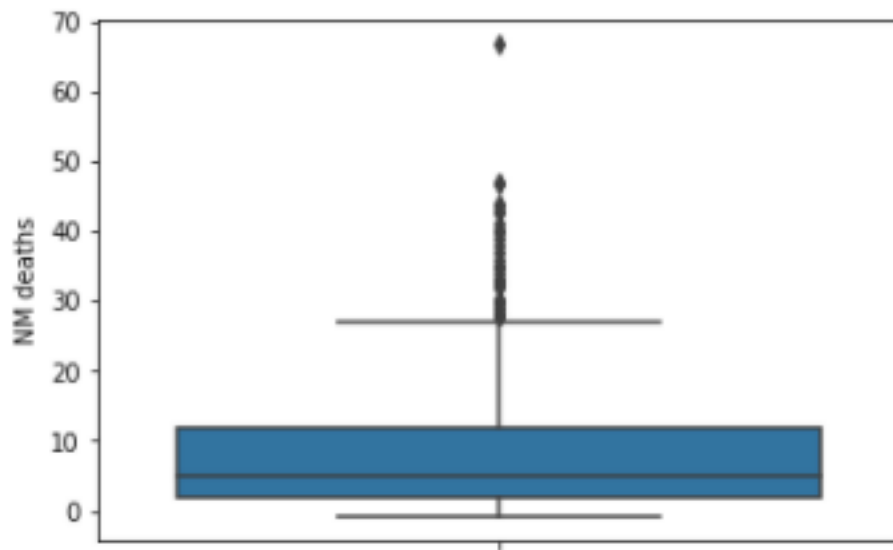
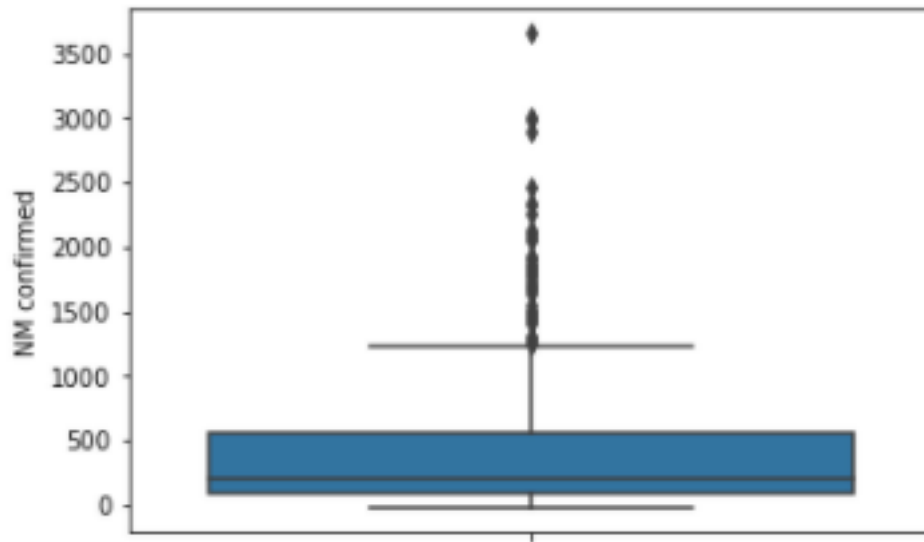
[438 rows x 5 columns]

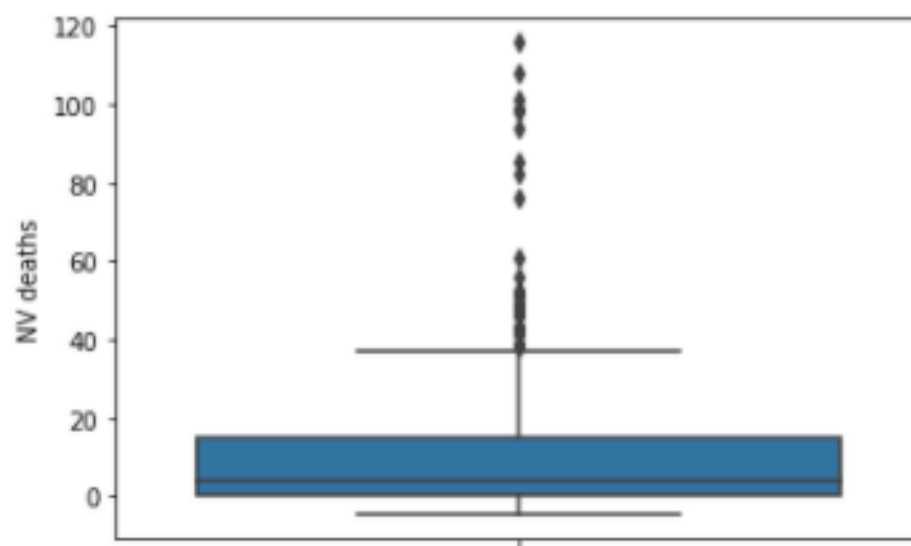
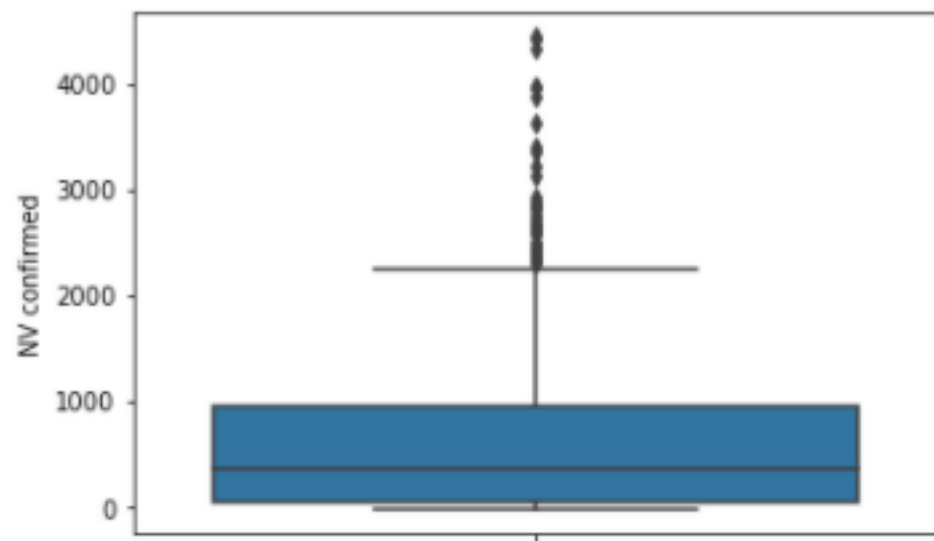
### Pre-Processed Data -

	Date	NM confirmed	NV confirmed	NM deaths	NV deaths
0	2020-01-22	0	0	0	0
1	2020-01-23	0	0	0	0
2	2020-01-24	0	0	0	0
3	2020-01-25	0	0	0	0
4	2020-01-26	0	0	0	0
..	...	...	...	...	...
433	2021-03-30	147	406	7	3
434	2021-03-31	275	276	5	10
435	2021-04-01	293	323	5	7
436	2021-04-02	208	0	7	0
437	2021-04-03	0	789	0	17

[438 rows x 5 columns]

## Box Plots





## Outlier Detection using Tukey's Rule -

	NM confirmed	NV confirmed	NM deaths	NV deaths
count	438.000000	438.000000	438.000000	438.000000
mean	438.712329	683.652968	9.015982	11.783105
std	587.542950	867.657454	10.603407	18.222022
min	-31.000000	-26.000000	-1.000000	-5.000000
25%	94.500000	42.250000	2.000000	0.000000
50%	197.000000	376.500000	5.000000	4.000000
75%	552.750000	948.000000	12.000000	15.000000
max	3665.000000	4455.000000	67.000000	116.000000

Outliers in  
NM confirmed  
are  
[289 290 292 293 294 295 299 300 301 302 303 304 305 306 307 308 309 310  
311 312 313 314 315 316 317 318 320 322 323 324 325 326 327 328 329 330  
331 332 335 337 338 343 344 345 350 351 352 353 358 359]

Outliers in  
NV confirmed  
are  
[302 306 307 308 310 311 314 316 317 318 319 320 321 322 325 326 327 329  
331 332 334 336 339 344 347 350 351 352 353 357 358 360 363 365 373]

Outliers in  
NM deaths  
are  
[300 305 307 310 313 315 316 317 318 322 324 326 328 329 330 331 336 337  
338 340 341 343 344 346 350 351 352 356 358 359 360 364 365 366 367 368  
372 378 385]

Outliers in  
NV deaths  
are  
[ 85 190 211 226 316 321 323 324 325 328 329 332 336 337 343 344 350 351  
352 353 357 358 360 365 366 367 370 371 373 374 378 379 380 382 386 394]

	NM confirmed	NV confirmed	NM deaths	NV deaths
count	438.000000	438.000000	438.000000	438.000000
mean	438.712329	683.652968	9.015982	11.783105
std	587.542950	867.657454	10.603407	18.222022
min	-31.000000	-26.000000	-1.000000	-5.000000
25%	94.500000	42.250000	2.000000	0.000000
50%	197.000000	376.500000	5.000000	4.000000
75%	552.750000	948.000000	12.000000	15.000000
max	3665.000000	4455.000000	67.000000	116.000000

## Removing Outliers –

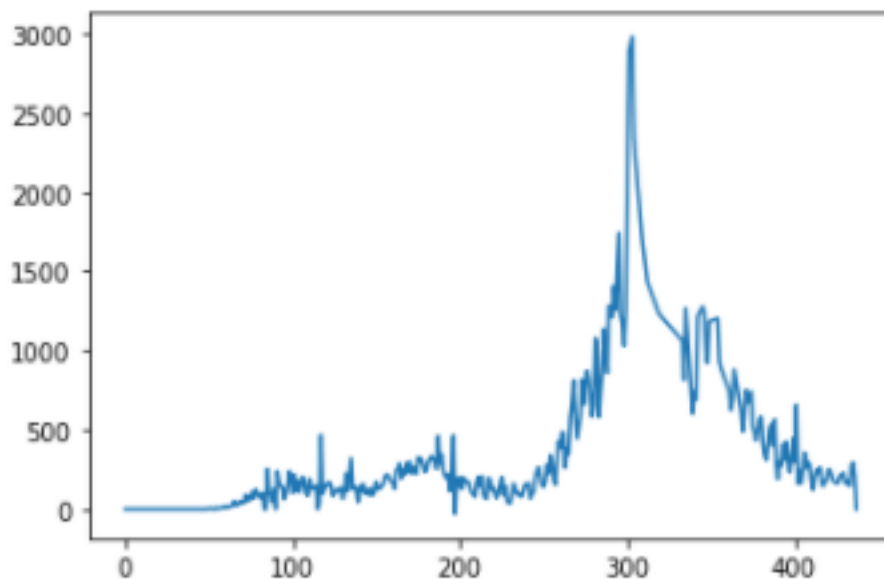
	NM confirmed	NV confirmed	NM deaths	NV deaths
count	438.000000	438.000000	438.000000	438.000000
mean	438.712329	683.652968	9.015982	11.783105
std	587.542950	867.657454	10.603407	18.222022
min	-31.000000	-26.000000	-1.000000	-5.000000
25%	94.500000	42.250000	2.000000	0.000000
50%	197.000000	376.500000	5.000000	4.000000
75%	552.750000	948.000000	12.000000	15.000000
max	3665.000000	4455.000000	67.000000	116.000000

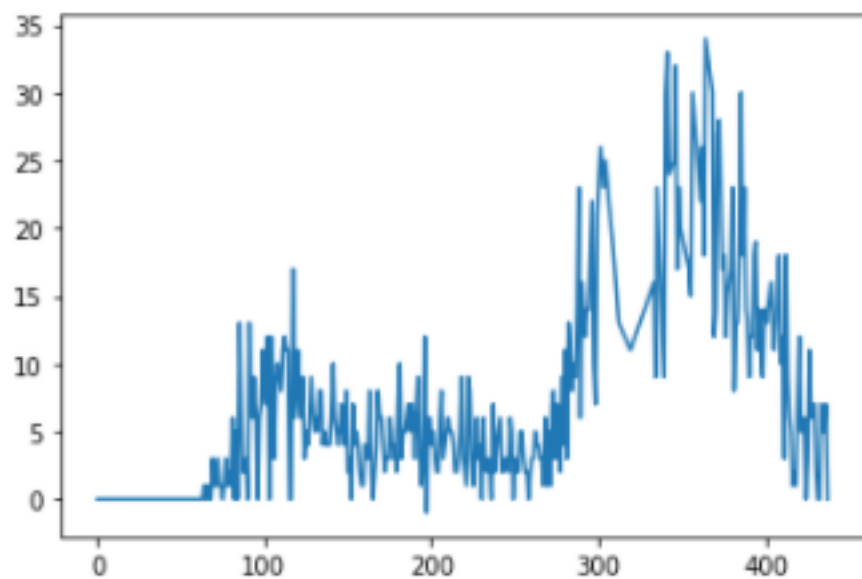
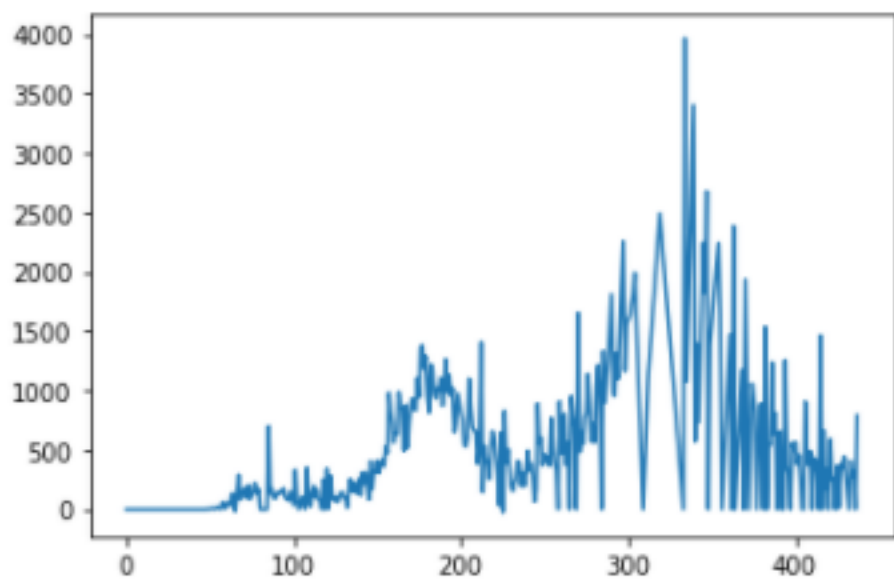
Dataset information after removing the outliers:

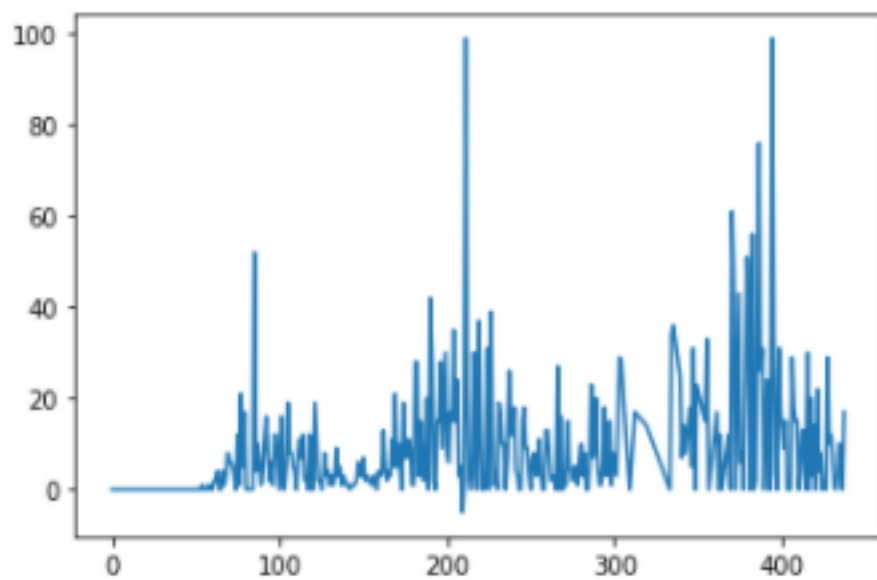
	NM confirmed	NV confirmed	NM deaths	NV deaths
count	393.000000	393.000000	393.000000	393.000000
mean	299.875318	478.783715	6.511450	8.511450
std	395.074666	563.658980	6.982522	12.948702
min	-31.000000	-26.000000	-1.000000	-5.000000
25%	83.000000	12.000000	1.000000	0.000000
50%	171.000000	318.000000	5.000000	3.000000
75%	315.000000	778.000000	9.000000	12.000000
max	2982.000000	3965.000000	34.000000	99.000000

As we can observe, initially there were 438 rows and now there are 393, we have removed 45 rows

## Plotting the Data









## Task 2A

Predicted values, MAPE and MSE for NM confirmed column using EWMA(0.5)

[23939.310861587524, 24120.655430793762, 24258.32771539688, 24363.66385769844, 24449.33192884922, 24590.66596442461, 24755.332982212305]

MSE: 78805.2932280279

MAPE: 1.1128022525383556

Predicted values, MAPE and MSE for NM confirmed column using EWMA(0.8)

[24056.59500545368, 24252.919001090744, 24367.383800218147, 24448.676760043625, 24517.735352008724, 24689.14707040175, 24873.829414080345]

MSE: 32279.585409043128

MAPE: 0.6896723334742848

Predicted values, MAPE and MSE for NV confirmed column using EWMA(0.5)

[63959.48520278931, 64514.24260139465, 65057.62130069733, 65533.81065034866, 65973.40532517433, 66319.70266258717, 66756.85133129358]

MSE: 933642.5217788888

MAPE: 1.4403479135973638

Predicted values, MAPE and MSE for NV confirmed column using EWMA(0.8)

[64618.69736554923, 64978.93947310982, 65476.587894622, 65903.3175789244, 66311.0635157849, 66595.01270315696, 67074.20254063139]

MSE: 317681.56909029145

MAPE: 0.8265661685839552

Predicted values, MAPE and MSE for NM deaths column using EWMA(0.5)

[733.9197020530701, 738.459851026535, 741.7299255132675, 744.3649627566338, 747.1824813783169, 751.0912406891584, 757.5456203445792]

MSE: 71.6982354634862

MAPE: 1.073503878707295

Predicted values, MAPE and MSE for NM deaths column using EWMA(0.8)

[737.7424089248121, 741.9484817849625, 744.3896963569923, 746.4779392713987, 749.2955878542797, 753.8591175708559, 761.9718235141711]

MSE: 30.976142479694722

MAPE: 0.6678169499320734

Predicted values, MAPE and MSE for NV deaths column using EWMA(0.5)

[1158.0162563323975, 1177.5081281661987, 1187.2540640830994, 1193.6270320415497, 1211.8135160207748, 1230.9067580103874, 1240.4533790051937]

MSE: 1053.4949132620761

MAPE: 2.4433640355936754

Predicted values, MAPE and MSE for NV deaths column using EWMA(0.8)

[1186.2458425590378, 1194.849168511808, 1196.5698337023614, 1199.3139667404719, 1223.8627933480946, 1244.7725586696188, 1248.9545117339235]

MSE: 461.6593640315824

MAPE: 1.332009556872688

Predicted values, MAPE and MSE for NM confirmed columns using AR

Predicted confirmed cases for NM confirmed with AR = 3: [24306.935545759785, 24514.85442225886, 24605.58961869839, 24669.361609357617, 24711.27165354359, 24914.875067647197, 25107.135634043312]

MSE: 7937.011344062339

MAPE: 0.2825468112371949

Predicted confirmed cases for NM confirmed with AR = 5: [24306.935545759785, 24514.85442225886, 24605.58961869839, 24669.361609357617, 24711.27165354359, 24914.875067647197, 25107.135634043312, 24359.32828098449, 24540.27493424469, 24642.05456086339, 24723.87331996068, 24775.506210868683, 24979.380284850347, 25182.139322956464]

MSE: 7937.011344062339

MAPE: 0.2825468112371949

Predicted values, MAPE and MSE for NV confirmed columns using AR

Predicted confirmed cases for NV confirmed with AR = 3: [65832.13649090889, 66110.32639109393, 66609.94469172033, 66832.26732795424, 67171.07726752959, 67327.81388016712, 67831.3441821913]

MSE: 235849.32362066637

MAPE: 0.6395136456177998

Predicted confirmed cases for NV confirmed with AR = 5: [65832.13649090889, 66110.32639109393, 66609.94469172033, 66832.26732795424, 67171.07726752959, 67327.81388016712, 67831.3441821913, 66368.3341437699, 66805.98688482025, 67328.72763242599, 67659.11311782707, 68038.39724203879, 67640.04242606706, 68089.9255220037]

MSE: 235849.32362066637

MAPE: 0.6395136456177998

Predicted values, MAPE and MSE for NM deaths columns using AR

Predicted confirmed cases for NM deaths with AR = 3: [747.6113449965665, 751.4746167337622, 752.6499780172844, 753.3930661948724, 755.584634924716, 760.6839342542276, 770.4726337130011]

MSE: 18.573941322172857

MAPE: 0.5232442616685796

Predicted confirmed cases for NM deaths with AR = 5: [747.6113449965665, 751.4746167337622, 752.6499780172844, 753.3930661948724, 755.584634924716, 760.6839342542276, 770.4726337130011, 748.2319077381833, 752.3652413120558, 753.8291875204968, 753.6929576193429, 755.1897961952396, 760.6956572778437, 771.8522291013041]

MSE: 18.573941322172857

MAPE: 0.5232442616685796

Predicted values, MAPE and MSE for NV deaths columns using AR

Predicted confirmed cases for NV deaths with AR = 3: [1226.9175390221262, 1228.5166375728802, 1221.164232169108, 1218.0219805160848, 1250.5288291304703, 1270.4395122750725, 1268.8081252368513]

MSE: 461.25041268512035

MAPE: 1.5658331016590075

Predicted confirmed cases for NV deaths with AR = 5: [1226.9175390221262, 1228.5166375728802, 1221.164232169108, 1218.0219805160848, 1250.5288291304703, 1270.4395122750725, 1268.8081252368513, 1258.3710650449075, 1259.8534962511285, 1260.2916823409864, 1257.7109268845436, 1282.9711912304465, 1295.3626772825141, 1288.4675623475077]

MSE: 461.25041268512035

MAPE: 1.5658331016590075

## Task 2B

	Date	NM confirmed	NV confirmed	NM deaths	NV deaths
376	2021-02-01	486	819	12	8
377	2021-02-02	432	0	15	0
379	2021-02-04	559	881	17	51
380	2021-02-05	582	888	23	39
381	2021-02-06	421	0	8	0
382	2021-02-07	343	1533	13	56
383	2021-02-08	311	0	13	0
384	2021-02-09	413	546	19	2
385	2021-02-10	509	541	30	36
386	2021-02-11	534	1236	18	76
387	2021-02-12	400	631	23	26
388	2021-02-13	565	810	14	31
389	2021-02-14	282	0	13	0
390	2021-02-15	190	514	9	15
391	2021-02-16	299	657	12	24
392	2021-02-17	272	0	12	0
393	2021-02-18	407	0	18	0
394	2021-02-19	312	1255	19	99
395	2021-02-20	425	363	11	34
396	2021-02-21	315	301	14	4
397	2021-02-22	233	0	11	0
398	2021-02-23	312	557	9	31
399	2021-02-24	445	506	14	17
400	2021-02-25	299	563	13	14
401	2021-02-26	657	380	14	9
402	2021-02-27	162	457	15	15
403	2021-02-28	240	265	16	0
	Date	NM confirmed	NV confirmed	NM deaths	NV deaths
404	2021-03-01	165	0	13	0
405	2021-03-02	244	0	11	0
406	2021-03-03	356	905	13	29
407	2021-03-04	259	376	16	18
408	2021-03-05	297	380	18	15
409	2021-03-06	282	488	10	15
410	2021-03-07	180	0	12	0
411	2021-03-08	121	429	3	4
412	2021-03-09	201	419	18	13
413	2021-03-10	247	318	9	13
414	2021-03-11	232	0	6	0
415	2021-03-12	262	1465	5	30
416	2021-03-13	182	0	1	0
417	2021-03-14	145	661	2	20
418	2021-03-15	178	179	1	3
419	2021-03-16	176	341	7	14
420	2021-03-17	243	0	12	0
421	2021-03-18	250	589	5	22
422	2021-03-19	218	256	5	4
423	2021-03-20	185	264	6	8

424	2021-03-21	171	0	0	0
425	2021-03-22	164	363	3	2
426	2021-03-23	167	11	11	0
427	2021-03-24	211	377	6	29
428	2021-03-25	217	345	7	10
429	2021-03-26	229	443	7	12
430	2021-03-27	169	297	2	8
431	2021-03-28	161	143	0	1
432	2021-03-29	182	0	0	0
433	2021-03-30	147	406	7	3
434	2021-03-31	275	276	5	10

## Wald's One Sample Test

```
NV confirmed
Reject NULL HYPOTHESIS 60.84468046051707
NM confirmed
Reject NULL HYPOTHESIS 67.27408494148042
NV deaths
Reject NULL HYPOTHESIS 23.24033583933362
NM deaths
Reject NULL HYPOTHESIS 16.413216173510026
```

### Hypothesis

Null hypothesis (H0): Mean of Feb 21 confirmed cases or deaths = Mean of March 21 confirmed cases or deaths.

Alternate hypothesis(H1): Mean of Feb 21 confirmed cases or deaths is not equal to mean of March 21 confirmed cases or deaths.

Procedure : We have taken the guess value as March 21 cases/deaths and  $\alpha = 0.05$  as given in documentation and sample mean as Feb 21. The standard error of the estimator is calculated in above walds function.

Result: W value for mean of Feb 21 NV confirmed cases =60.84468046051707 which is greater than 1.96 we are rejecting the NULL hypothesis. W value for mean of Feb 21 NM confirmed cases =67.27408494148042 which is greater than 1.96 we are rejecting the NULL hypothesis. W value for mean of Feb 21 NV deaths =23.24033583933362 which is greater than 1.96 we are rejecting the NULL hypothesis. W value for mean of Feb 21 NM deaths =16.413216173510026 which is greater than 1.96 we are rejecting the NULL hypothesis.

Is Test Applicable ? The main Assumptions of Wald's test is that the sample data has to be normally distributed. Since we are using a mean estimator which is Poisson MLE, using CLT we can say that the data is asymptotically normal.

Hence ,We can conclude the Wald's Test is applicable on given dataset.

## Walds Two Sample Test

```
NV confirmed
Reject NULL HYPOTHESIS 36.00135062635395
NM confirmed
Reject NULL HYPOTHESIS 38.17819402331588
NM deaths
Reject NULL HYPOTHESIS 12.026445802057546
NV deaths
Reject NULL HYPOTHESIS 8.88072270686856
```

### Hypothesis

Null hypothesis (H0): Mean of Feb 21 confirmed cases or deaths = Mean of March 21 confirmed cases or deaths.

Alternate hypothesis(H1): Mean of Feb 21 confirmed cases or deaths is not equal to mean of March 21 confirmed cases or deaths.

Procedure : We have taken the  $\alpha = 0.05$  as given in documentation and calculated the numerator and denominator of  $w$  in the above `walds_2_sample_testing` function . The standard error of the estimator is combination of the standard error of both the months data which is February 21 and March 21.

Result: W value for mean of Feb 21 NV confirmed cases =36.00135062635395 which is greater than 1.96 we are rejecting the NULL hypothesis. W value for mean of Feb 21 NM confirmed cases =38.17819402331588 which is greater than 1.96 we are rejecting the NULL hypothesis. W value for mean of Feb 21 NM deaths =12.026445802057546 which is greater than 1.96 we are rejecting the NULL hypothesis. W value for mean of Feb 21 NV deaths =8.88072270686856 which is greater than 1.96 we are rejecting the NULL hypothesis.

Is Test Applicable ? The main Assumptions of Wald's test is that the sample data has to be normally distributed. Since we are using a mean estimator, using CLT we can say that the data is asymptotically normal.

Hence ,We can conclude the Wald's Test is applicable on given dataset.

## Z Test

```
NV confirmed
Accept NULL HYPOTHESIS 0.011833444930198568
NM confirmed
Accept NULL HYPOTHESIS 0.01463507248320475
NV deaths
Accept NULL HYPOTHESIS 0.049608858754022205
NM deaths
Accept NULL HYPOTHESIS 0.036355495694683625
```

## Hypothesis

Null hypothesis (H0): Mean of Feb 21 confirmed cases or deaths = Mean of March 21 confirmed cases or deaths.

Alternate hypothesis(H1): Mean of Feb 21 confirmed cases or deaths is not equal to mean of March 21 confirmed cases or deaths.

Procedure : We have taken the guess value as March 21 cases/deaths and alpha = 0.05 as given in documentation and sample mean as Feb 21. We used the corrected sample standard deviation of the entire COVID19 dataset we had for each state as the true sigma value.

Result: Z value for mean of Feb 21 NV confirmed cases = 0.011833444930198568 which is less than 1.96 we are accepting the NULL hypothesis. Z value for mean of Feb 21 NM confirmed cases = 0.01463507248320475 which is less than 1.96 we are accepting the NULL hypothesis. Z value for mean of Feb 21 NM deaths = 0.036355495694683625 which is less than 1.96 we are accepting the NULL hypothesis. Z value for mean of Feb 21 NV deaths = 0.049608858754022205 which is less than 1.96 we are accepting the NULL hypothesis.

Is Test Applicable ? The main Assumptions of Z-test are that true standard deviation is known to us, sample size has to be large or the sample data has to be normally distributed. The true standard deviation is known to us. The data here has 400+ rows so it is large. Since we are using a mean estimator, using CLT we can say that the data is asymptotically normal.

Hence ,We can conclude the Z-Test is applicable on given dataset.



## T Test

NV confirmed  
Reject NULL HYPOTHESIS 3.466232317733068  
NM confirmed  
Reject NULL HYPOTHESIS 18.584110592851285  
NV deaths  
Reject NULL HYPOTHESIS 7.377885488184716  
NM deaths  
Reject NULL HYPOTHESIS 8.436562085344791

## Hypothesis

Null hypothesis (H0): Mean of Feb 21 confirmed cases or deaths = Mean of March 21 confirmed cases or deaths.

Alternate hypothesis(H1): Mean of Feb 21 confirmed cases or deaths is not equal to mean of March 21 confirmed cases or deaths.

Procedure : We have taken the  $\alpha = 0.05$  as given in documentation and degree of freedom as 30 to calculate the value of T.

Result: T value for mean of Feb 21 NV confirmed cases = 3.466232317733068 which is greater than 1.697261 we are rejecting the NULL hypothesis. T value for mean of Feb 21 NM confirmed cases = 18.584110592851285 which is greater than 1.697261 we are rejecting the NULL hypothesis. T value for mean of Feb 21 NM deaths = 8.436562085344791 which is greater than 1.697261 we are rejecting the NULL hypothesis. T value for mean of Feb 21 NV deaths = 7.377885488184716 which is greater than 1.697261 we are rejecting the NULL hypothesis.

Is Test Applicable ? The T-test is not valid since the data points are expected to follow a Normal distribution but the given distribution to us is Poisson.

## Unpaired T Test

NV confirmed  
Accept NULL HYPOTHESIS 1.4149094874966222  
NM confirmed  
Reject NULL HYPOTHESIS 5.141872121508538  
NV deaths  
Accept NULL HYPOTHESIS 1.9104971125134882  
NM deaths  
Reject NULL HYPOTHESIS 4.223523081107441

## Hypothesis

Null hypothesis (H0): Mean of Feb 21 confirmed cases or deaths = Mean of March 21 confirmed cases or deaths.

Alternate hypothesis(H1): Mean of Feb 21 confirmed cases or deaths is not equal to mean of March 21 confirmed cases or deaths.

Procedure : We have taken the  $\alpha = 0.05$  as given in documentation and degree of freedom as 30 to calculate the value of T.

Result: T value for mean of Feb 21 NV confirmed cases = 1.4149094874966222 which is less than 2.000995 we are accepting the NULL hypothesis. T value for mean of Feb 21 NM confirmed cases = 5.141872121508538 which is greater than 2.000995 we are rejecting the NULL hypothesis. T value for mean of Feb 21 NM deaths = 4.223523081107441 which is greater than 2.000995 we are rejecting the NULL hypothesis. T value for mean of Feb 21 NV deaths = 1.9104971125134882 which is less than 2.000995 we are accepting the NULL hypothesis.

Is Test Applicable ? The T-test however is not valid since the data points are expected to follow a Normal distribution but the given distribution to us is Poisson.

## Task 2C

### KS test

Checking equality of distributions for confirmed cases in 2 states using Poisson distribution  
mme\_lambda: 916.5

Maximum Difference: 0.46659659496804107  
Null hypothesis is rejected as Oct-Dec 2020 data for the second state does not have the distribution with the obtained MME parameters for Confirmed cases

Checking equality of distributions for confirmed cases in 2 states using Geometric distribution  
mme\_p: 0.0010911074740861974

Maximum Difference: 0.23595778230524775  
Null hypothesis is rejected as Oct-Dec 2020 data for the second state does not have the distribution with the obtained MME parameters for Confirmed cases

Checking equality of distributions for confirmed cases in 2 states using Binomial distribution  
mme\_p: -354.76375704673575  
mme\_n: -2.583409330280777

Maximum Difference: 1.0  
Null hypothesis is rejected as Oct-Dec 2020 data for the second state does not have the distribution with the obtained MME parameters for Confirmed cases

Checking equality of distributions for deaths in 2 states using Poisson distribution  
mme\_lambda: 17.391304347826086

Maximum Difference: 0.4700505411726082  
Null hypothesis is rejected as Oct-Dec 2020 data for the second state does not have the distribution with the obtained MME parameters for Deaths

Checking equality of distributions for deaths in 2 states using Geometric distribution  
mme\_p: 0.0575

Maximum Difference: 0.11956521739130432  
Null hypothesis is rejected as Oct-Dec 2020 data for the second state does not have the distribution with the obtained MME parameters for Deaths

Checking equality of distributions for deaths in 2 states using Binomial distribution  
mme\_p: -10.008695652173913  
mme\_n: -1.7376194613379672

Maximum Difference: 1.0

Null hypothesis is rejected as Oct-Dec 2020 data for the second state does not have the distribution with the obtained MME parameters for Deaths

## Permutation Test

Permutation test for daily confirmed data for New Mexico and Nevada

observed\_T= 110.96666666666667

alpha = 0.05

For n = 1000 random permutations, p\_value: 0.374

Therefore, NULL hypothesis for 1000 permutations can be rejected as p-value is less than alpha

Permutation test for daily deaths data for Georgia and Hawaii

observed\_T= 1.3152173913043477

alpha = 0.05

For n = 1000 random permutations, p\_value: 0.565

Therefore, NULL hypothesis for 1000 permutations can be rejected as p-value is less than alpha

## Task 2D

Lambda after first 4 weeks is 7.392857142857143  
Lambda after first 5 weeks is 7.485714285714286  
Lambda after first 6 weeks is 8.404761904761905  
Lambda after first 7 weeks is 8.83673469387755  
Lambda after first 8 weeks is 10.053571428571429

