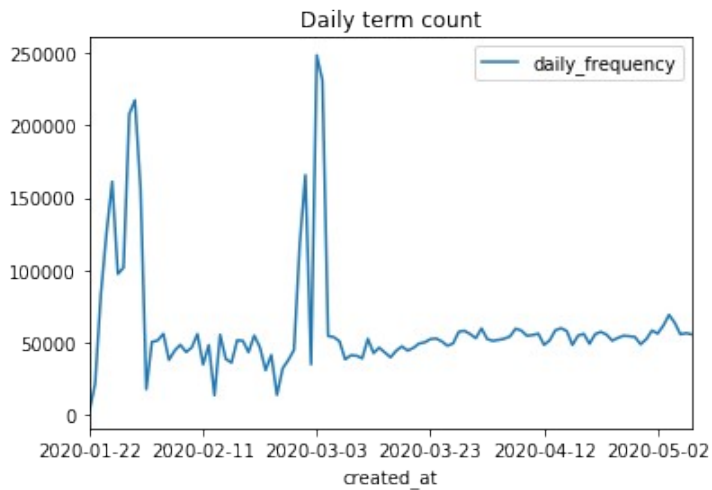


## Outbreak and Symptom Term analysis

The data used is a subset of the original data from 22-01-2020 to 08-05-2020 and consists of ~11 million rows.



### 1) Daily Count

Unusually large counts (>100000) were observed on -

3	2020-01-25
4	2020-01-26
6	2020-01-28
7	2020-01-29
8	2020-01-30
9	2020-01-31
37	2020-02-29
38	2020-03-01
40	2020-03-03
41	2020-03-04

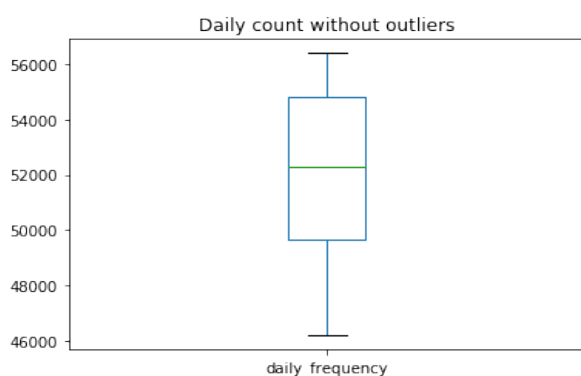
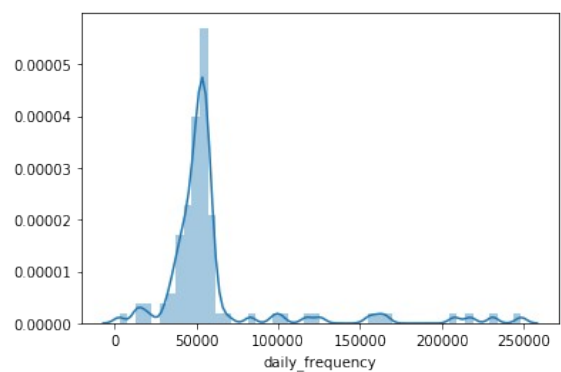
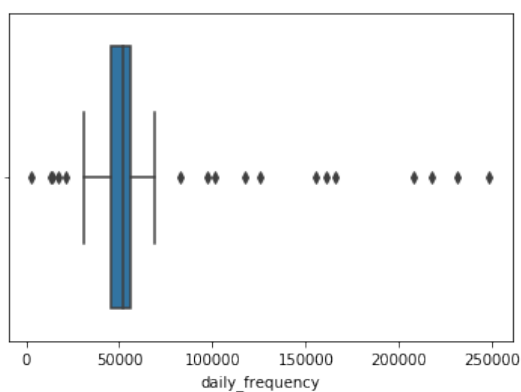
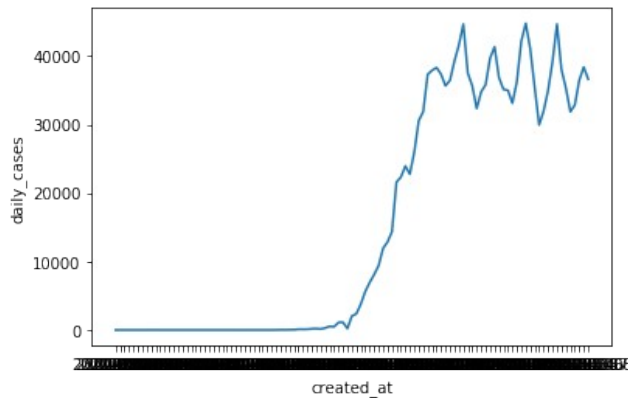


Fig. Distribution plot of daily term count

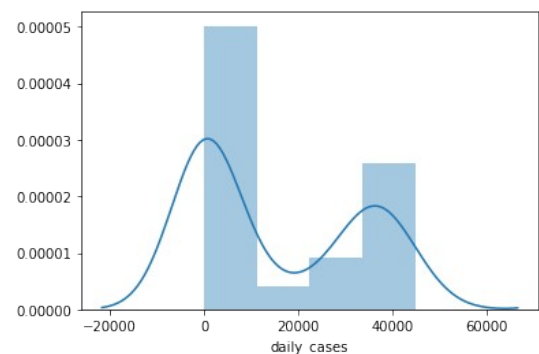
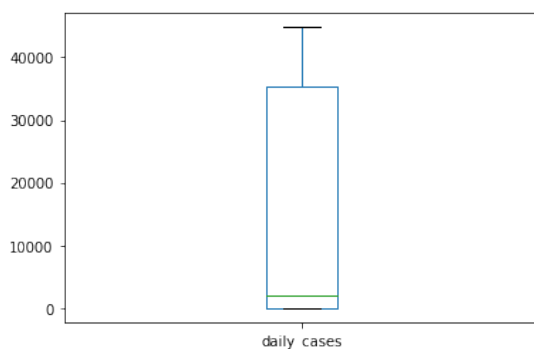
From the boxplot and distribution plot, median count is nearly 52000 and distribution is more or less symmetric if outliers are ignored.

Note : “Daily cases” and “daily deaths” were obtained from <https://datahub.io/core/covid-19#readme>. Daily deaths and daily cases were obtained for the countries - 'Australia', 'Canada', 'India', 'New Zealand', 'US', 'United Kingdom'

## 2) Infection Rate



Initially, number of cases reported were low, but as testing ramped up, daily\_cases began rising too.



No outlier is apparent from the boxplot, but the distribution is highly skewed with median around 2000.

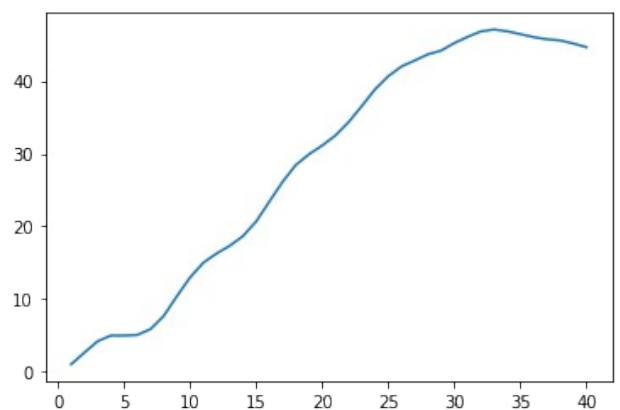
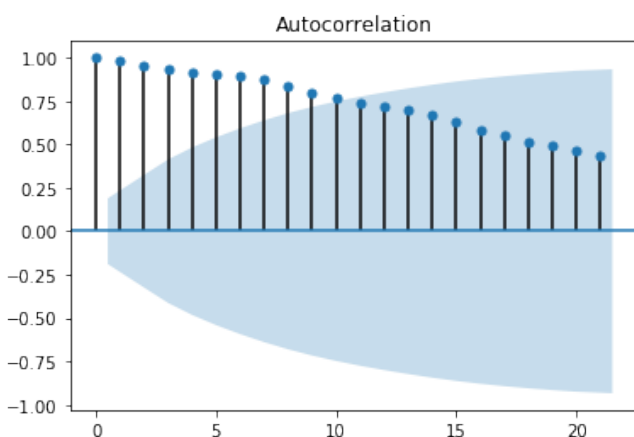
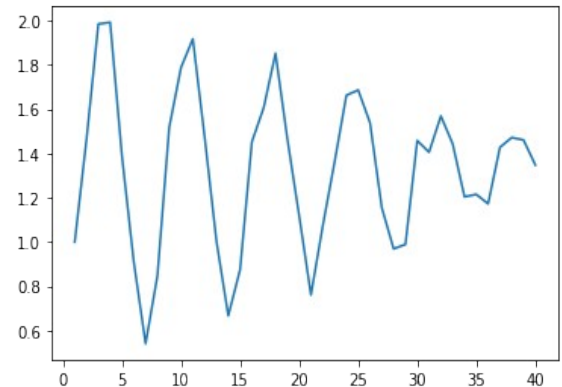
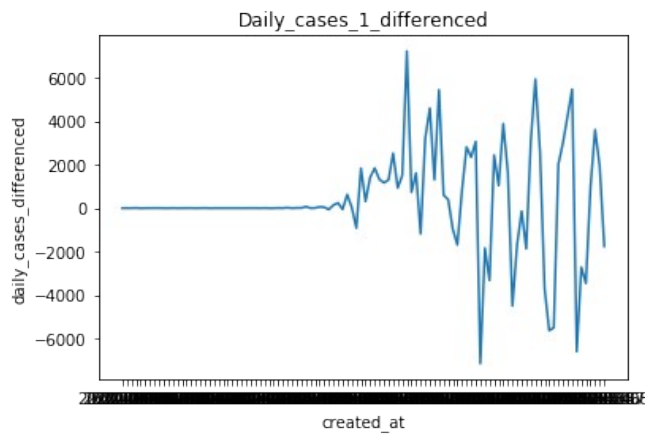


Fig – Variogram analysis

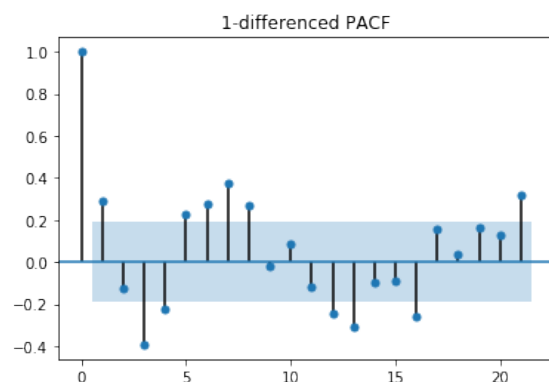
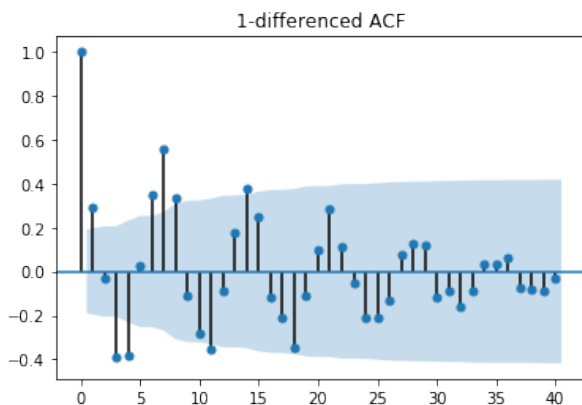
In the ACF plot, the acf values are significant and cutoff at high lags. Also the variogram values also increase with lags. This suggests that the “infection\_rate” is nonstationary and needs at least one order of differencing.

The differenced plot looks like -



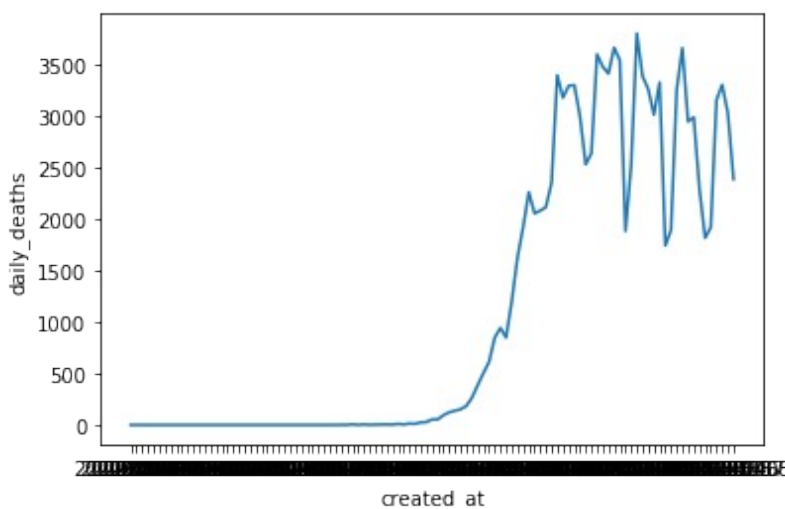
Variogram analysis of differenced series

The differenced series looks stationary with mean around 0.

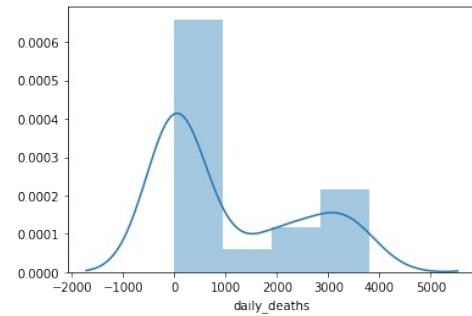
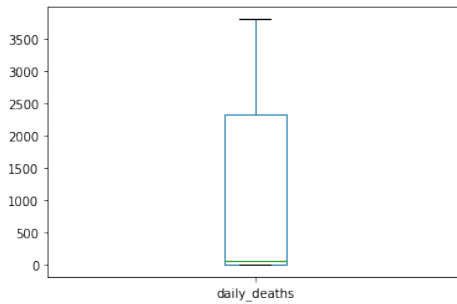


Both ACF and PACF plots show a damped sine wave which suggests an ARMA process of suitable order. Thus, daily\_cases may be modelled by ARIMA process.

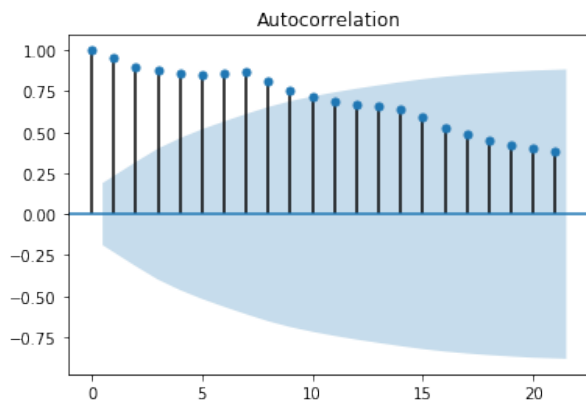
## 2) Daily Deaths



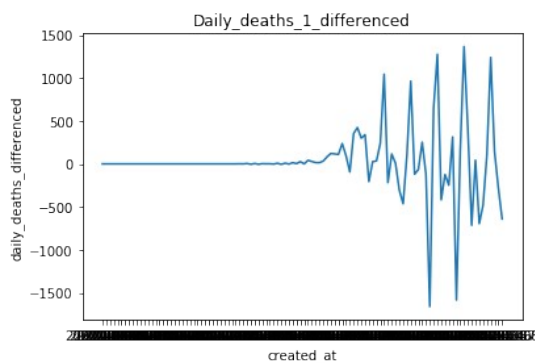
As expected, daily\_deaths rise with daily cases.



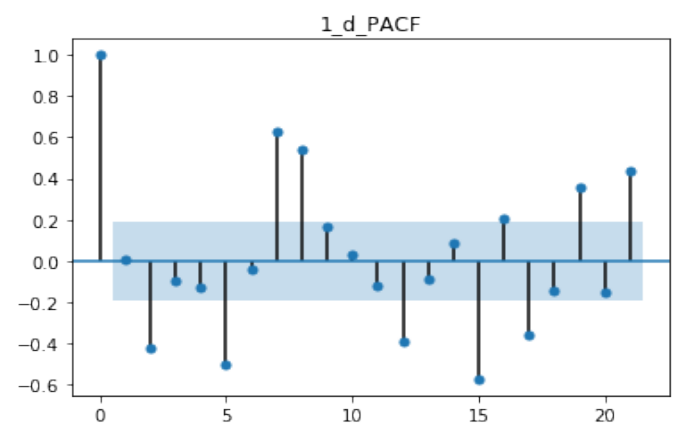
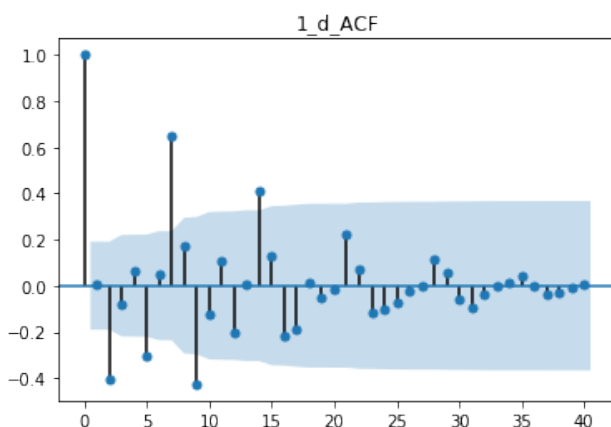
From the boxplot, no apparent outlier is observed. However, the distribution of daily deaths is highly skewed with median around  $<100$  as seen in distribution plot as well as boxplot.



The acf plots show significant acf values at high lags which suggests non-stationarity in time series.

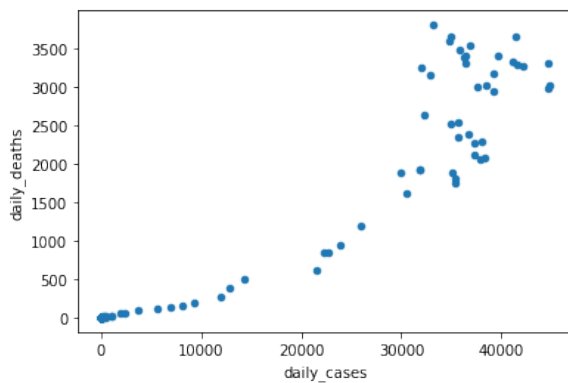
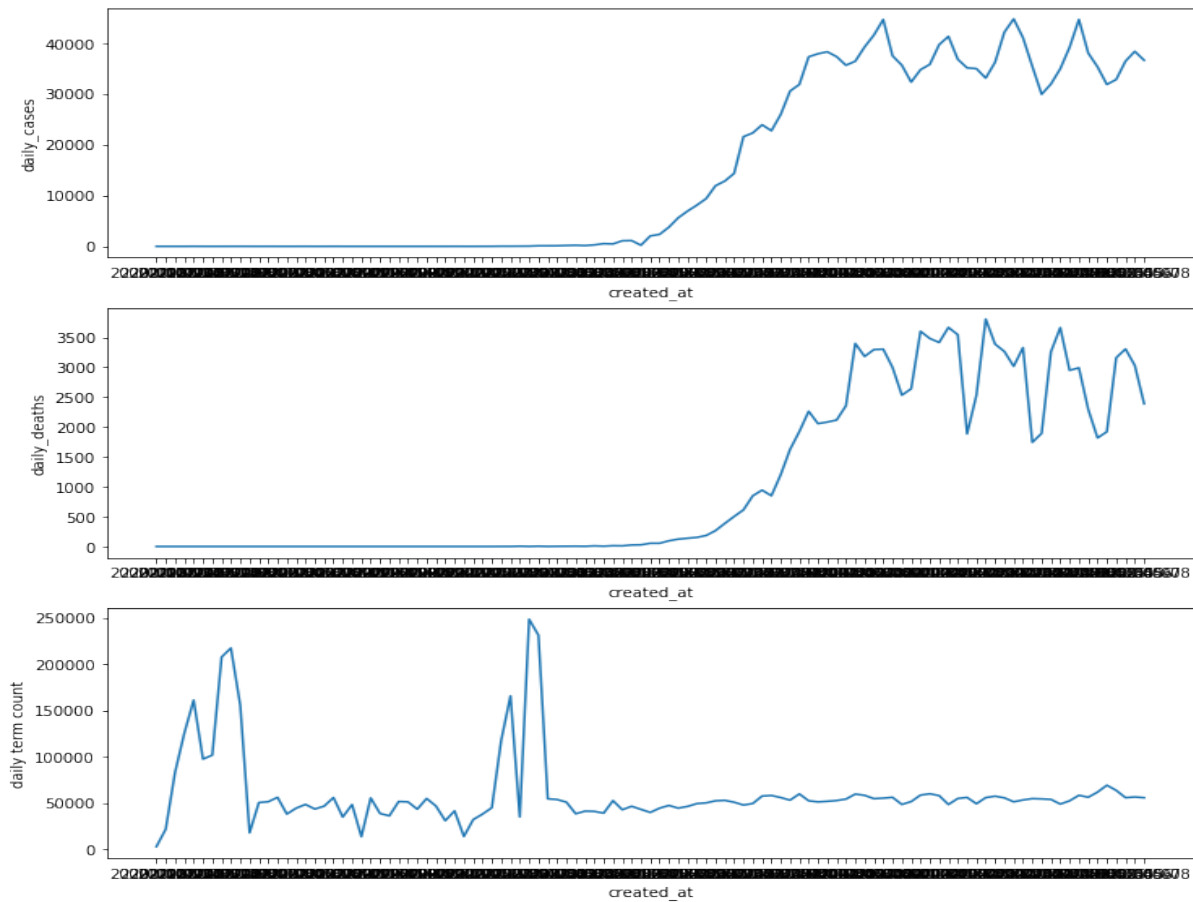


After differencing, daily deaths resembles a stationary time series.



From ACF and PACF plots, it seems the differenced daily\_deaths time series follows ARMA process of suitable order.

Time series plots of daily\_deaths,daily\_cases and daily term count are plotted below



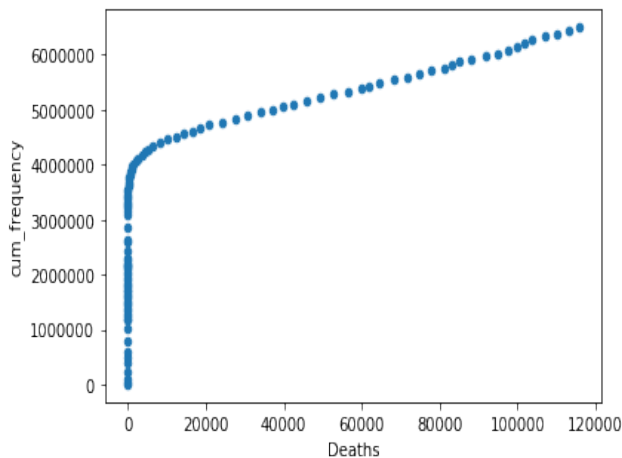
From the above graphs, daily\_death and daily\_cases look similar and have high correlation coefficient 0.9509412732821179, which is expected because as daily cases rise, daily deaths increase too.

However, daily\_count doesn't have much correlation with daily\_cases and daily\_deaths with the respective correlation coefficients being -

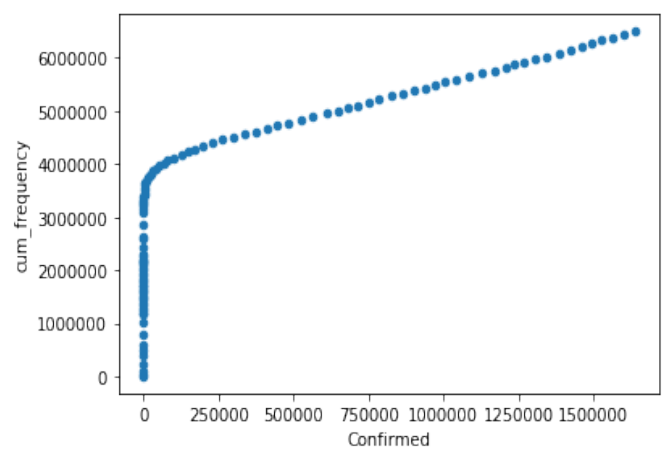
-0.12905853084101224

-0.10939421802504366

But, there is high correlation between “cumulative term counts” and “cumulative cases” and “cumulative deaths” as shown below



correlation-coefficient : 0.7973005485418727

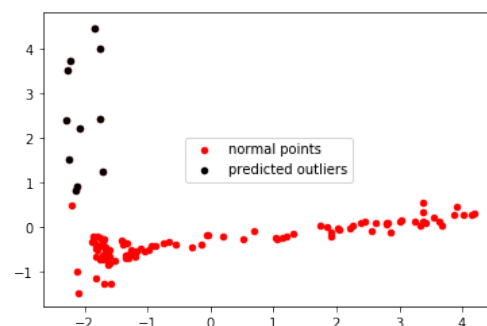


correlation-coefficient : 0.8281780247336641

## Anomaly Detection

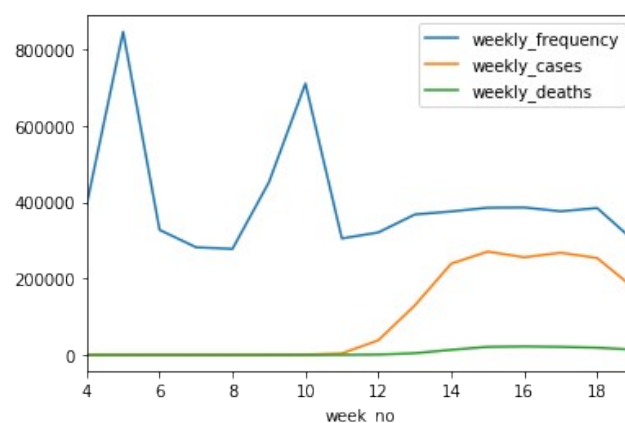
Using the features- daily\_cases, daily\_term\_count, daily\_deaths, cumulative\_term\_count, cumulative\_deaths, cumulative\_cases

sklearn's Local Outlier Factor algorithm was used on normalised data. The features were projected on 2-D space using PCA.



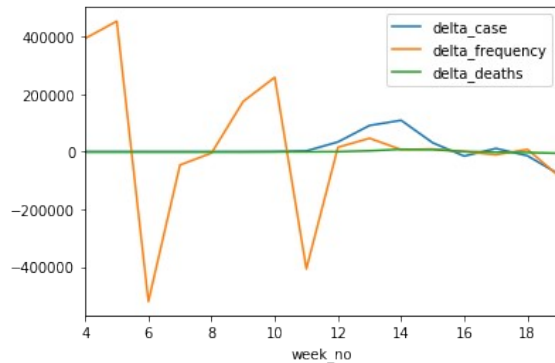
Those points were predicted as outliers for which daily\_cases and daily\_deaths were low, but the daily\_term\_counts were high.

## Weekly analysis



Weekly data represents a similar trend as daily data. Weekly count of terms have unusually high peaks at week no.s 5 and 10, but is below 400000 in the remaining weeks.

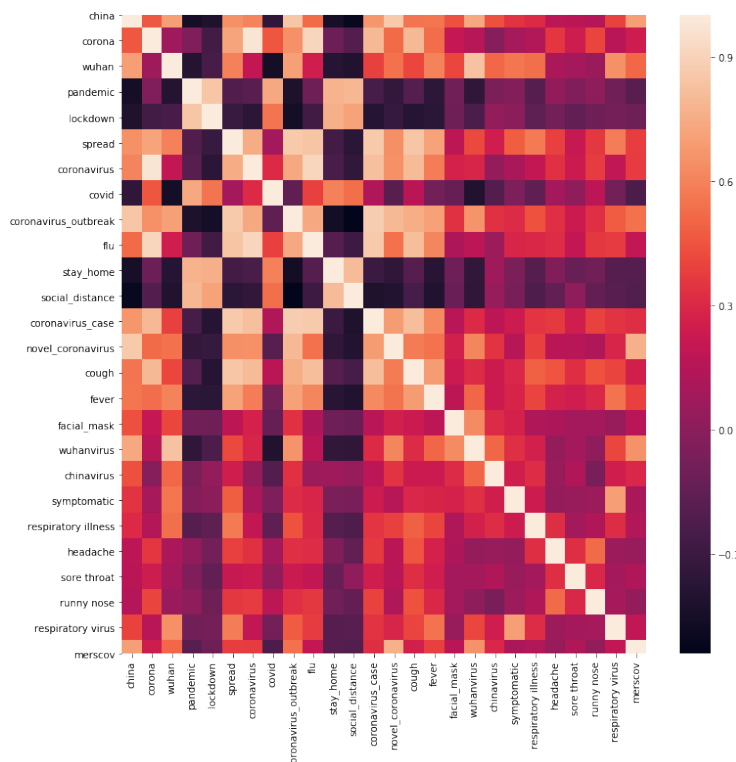
Weekly cases show a high correlation with weekly\_deaths with correlation coefficient being 0.9754733321404416, whereas weekly\_cases and weekly\_deaths continue to have low correlation with weekly term count.



Weekly change in term counts fluctuate between positive and negative values, whereas weekly change in deaths and cases is nearly zero across the weeks. However, there is a surge in weekly cases from week 12 to week 14.

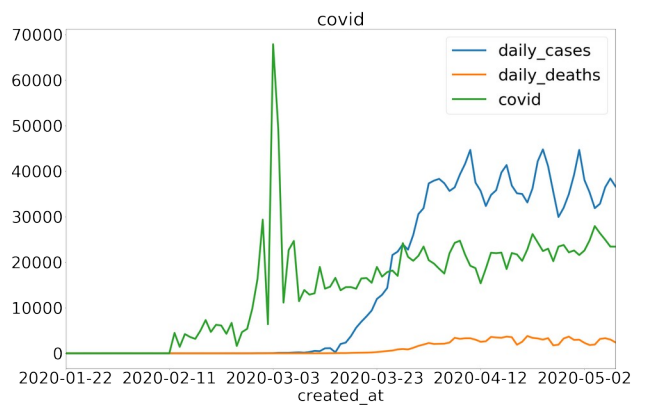
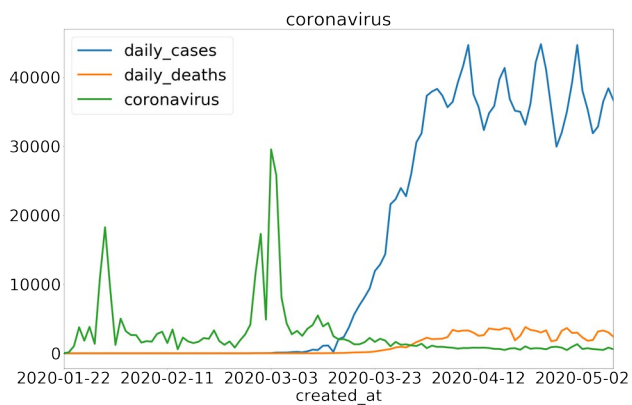
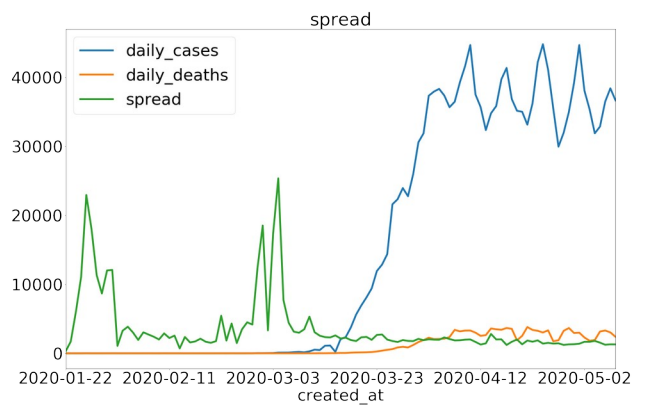
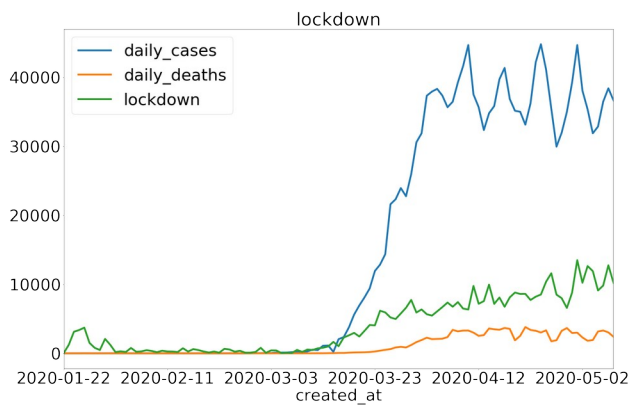
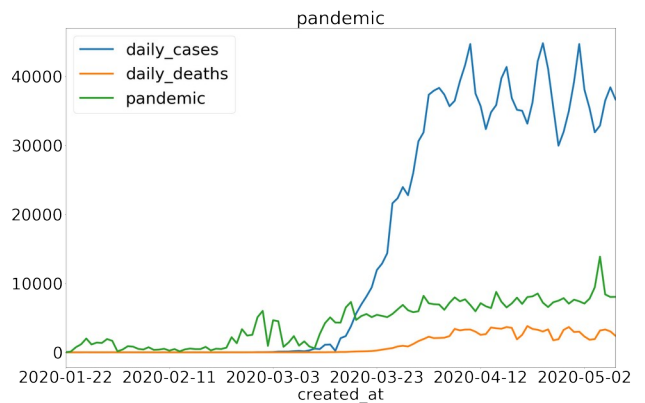
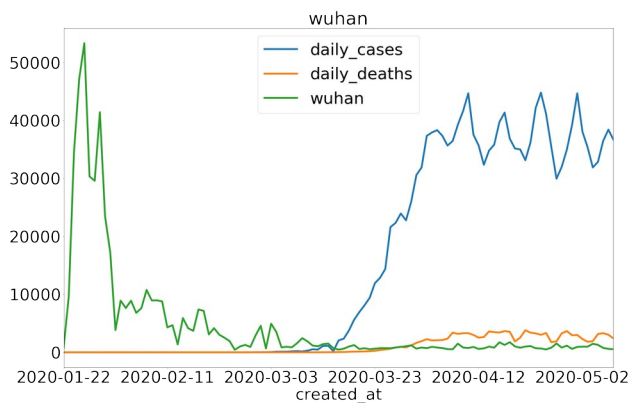
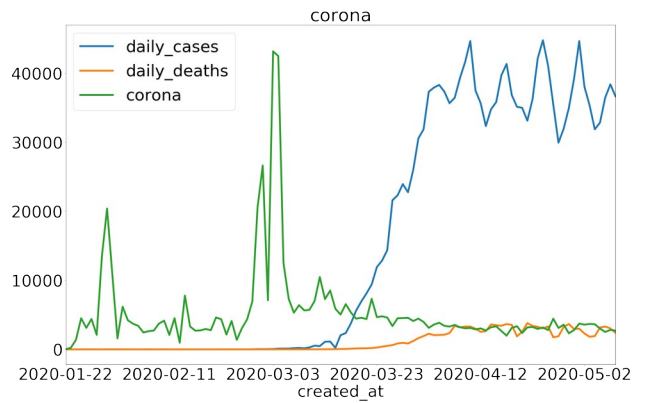
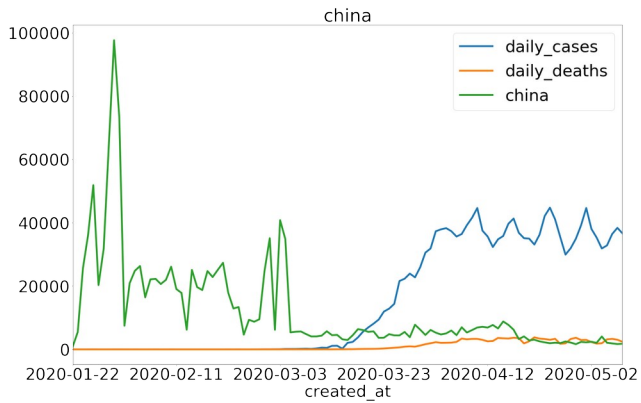
## Termwise-Analysis

A set of outbreak and symptom terms were identified and the frequency of occurrences of those terms in individual tweets were obtained. From the set of terms, top frequently occurring terms were chosen on which further analysis has been done.

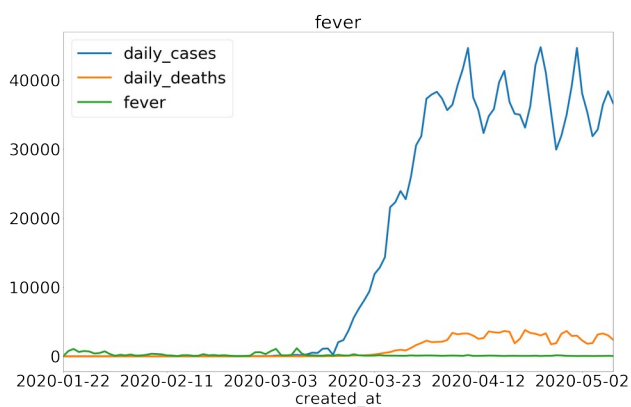
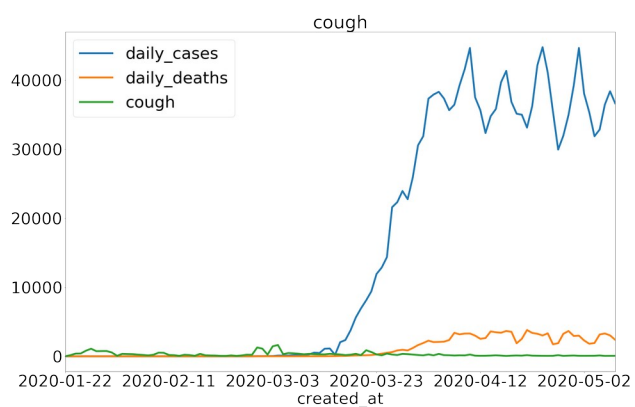
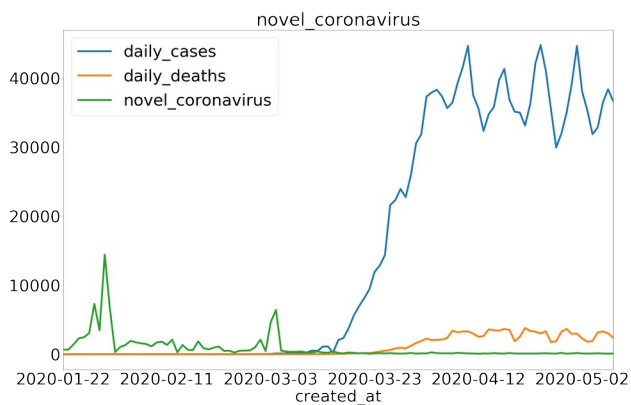
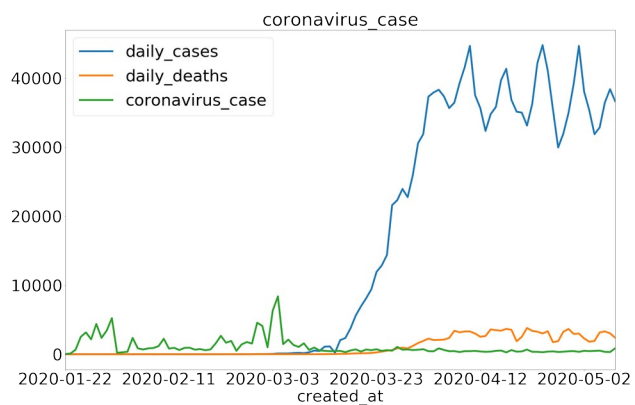
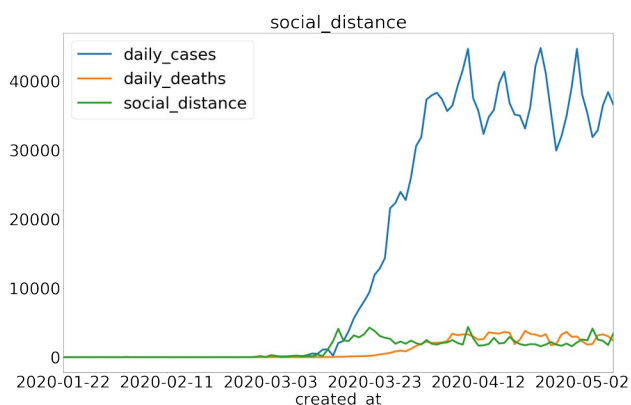
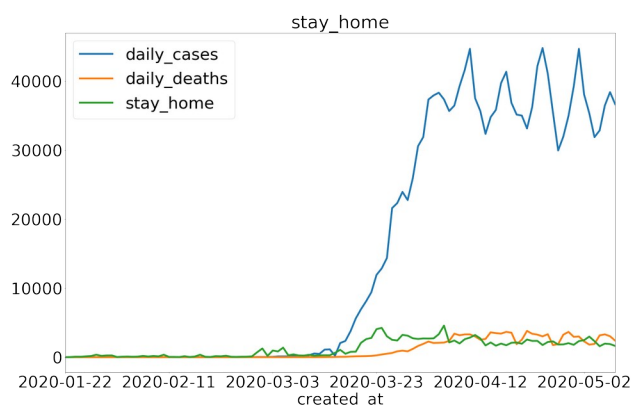
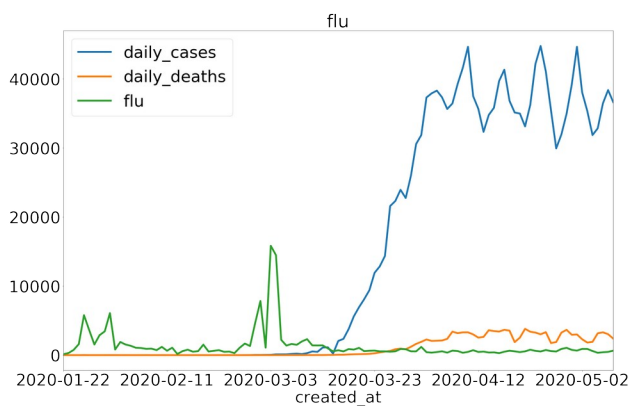
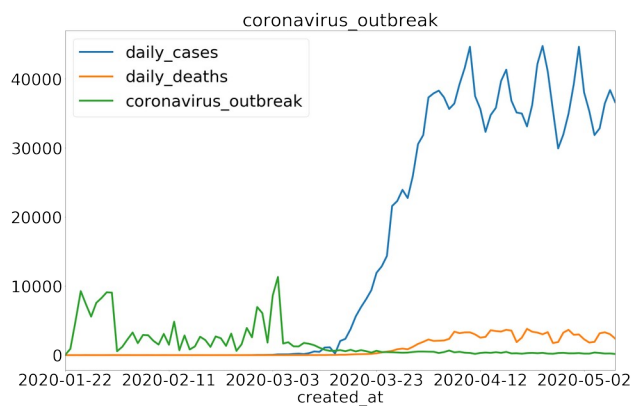


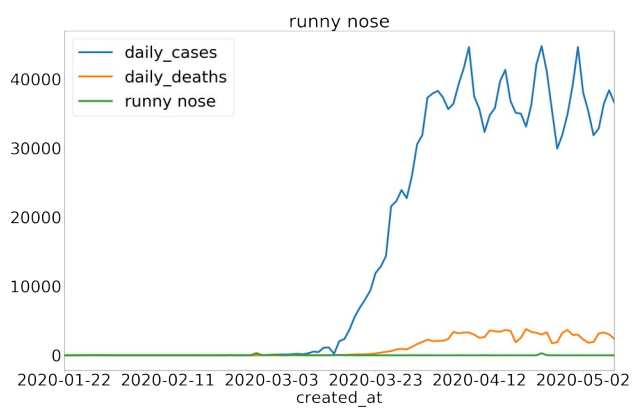
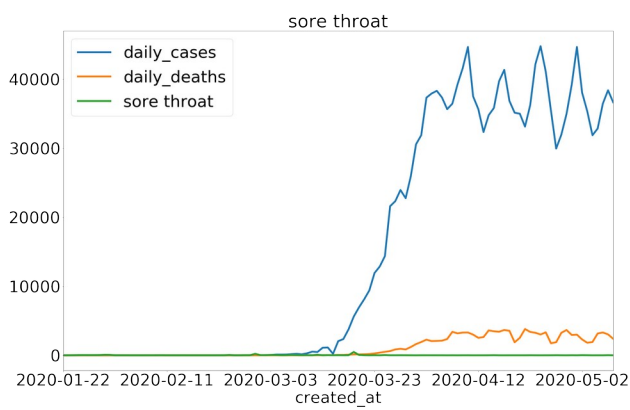
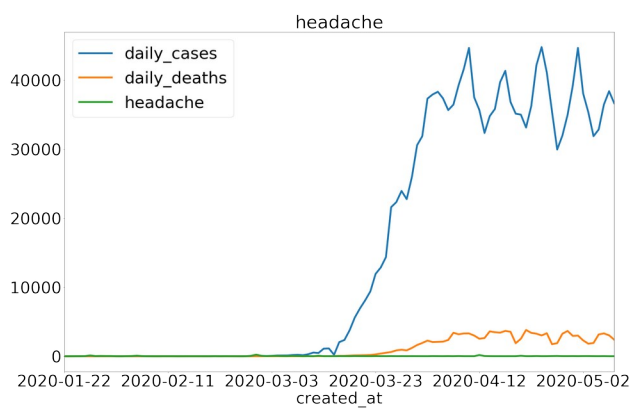
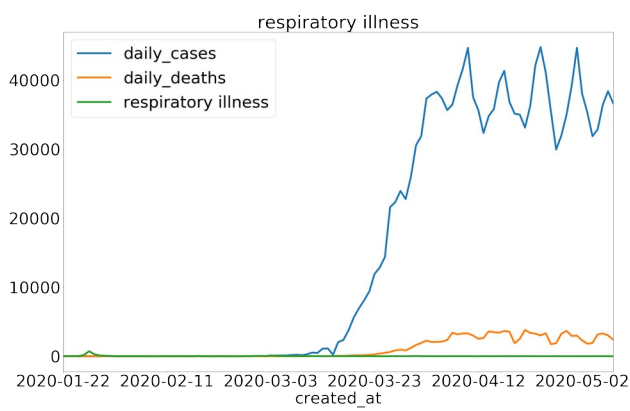
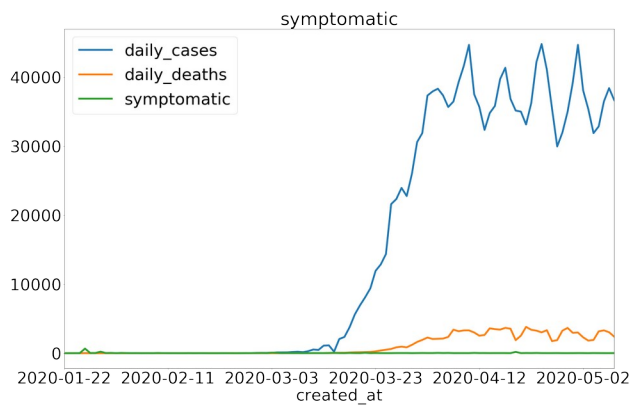
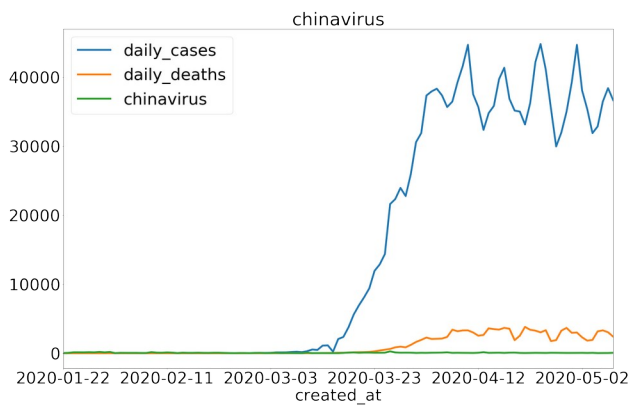
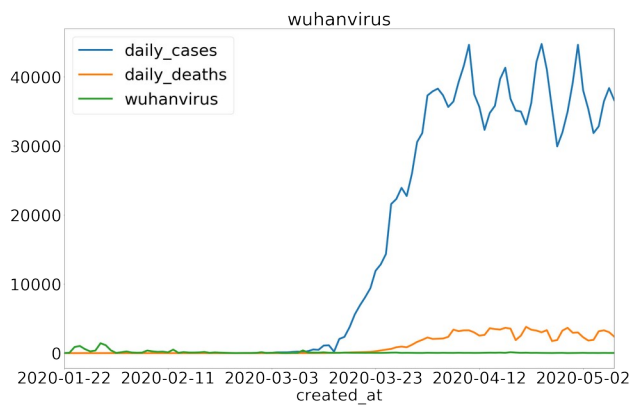
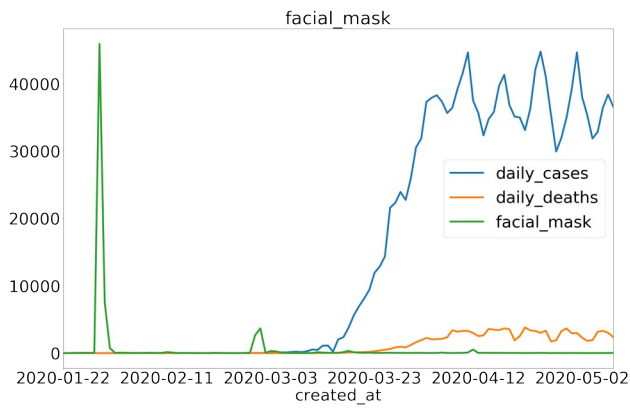
It can be seen that some terms have high correlation with other terms whereas some have little or no correlation at all.

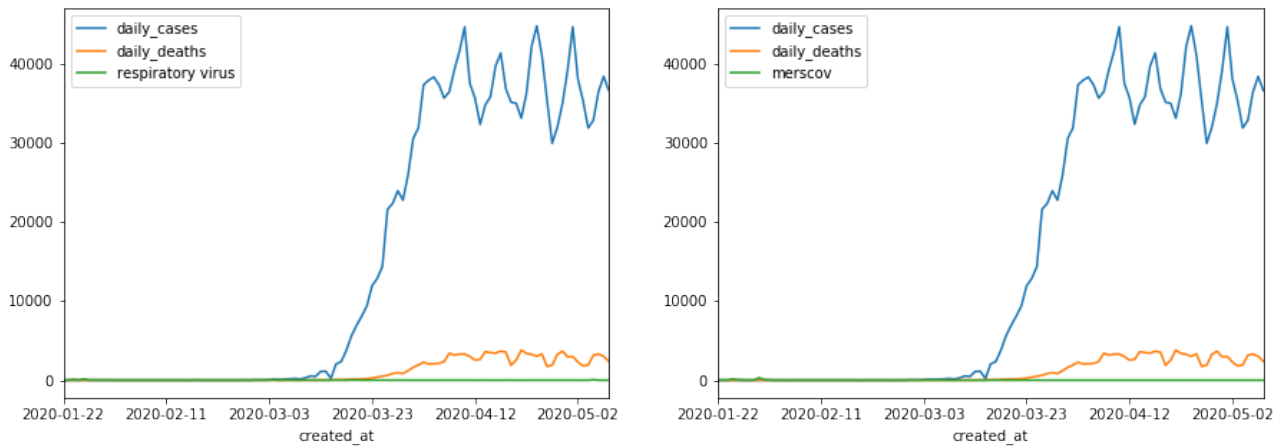
The time series plots below is plotted for each of those chosen terms:











From the time series plots above, it can be seen that certain keywords like “china”, “corona” were highly used in tweets before cases and deaths started being reported. But, as the number of cases began rising, their usage started declining.

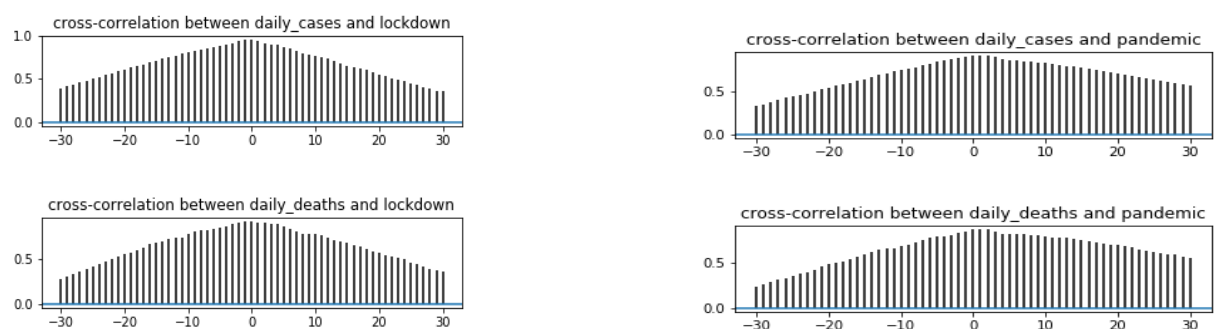
However, usage of keywords like “stay\_home”, “social\_distance” started picking up as cases and deaths started rising.

This pattern in frequency of keywords therefore reflects the perception of disease in conjunction with the evolution of the outbreak.

To see the correlation of the individual time series pattern of words with infections and deaths, cross-correlation plot was plotted for each term to identify the lags at which they correlate the most with the targets “infection” and “death”.

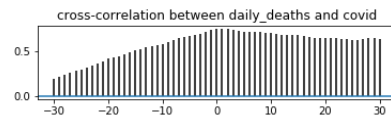
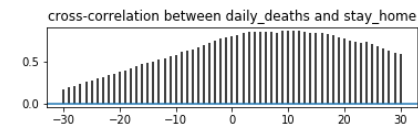
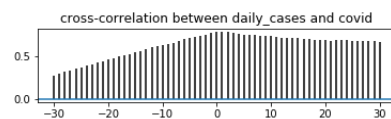
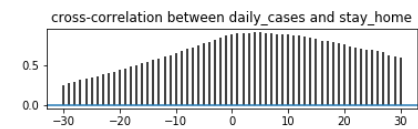
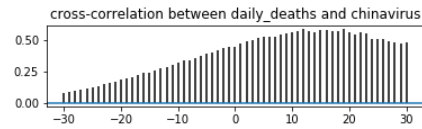
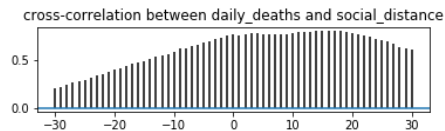
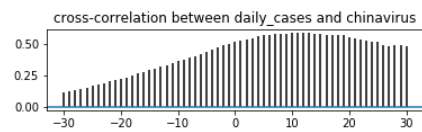
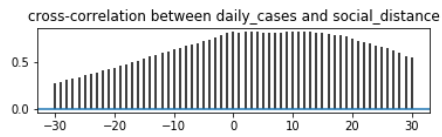
Terms showing similar pattern and range of values were then clubbed together as single term. Altogether 8 such “clubbed\_terms” were obtained:

### 1)Term 1

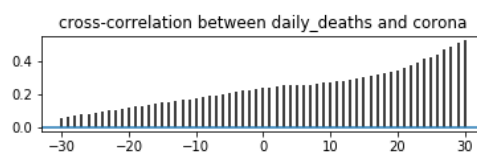
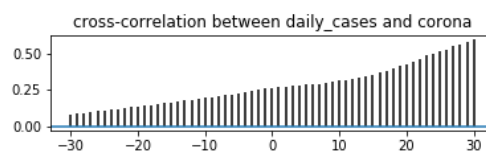
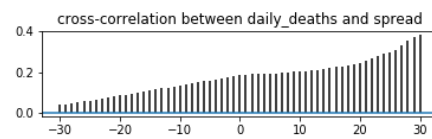
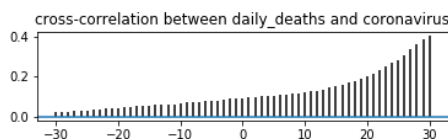
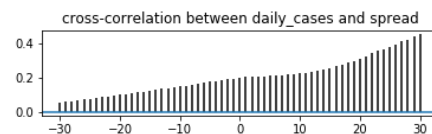
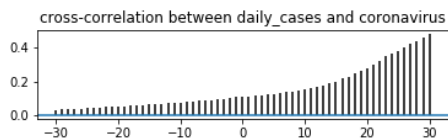


“pandemic” and “lockdown” show similar pattern. The correlation is maximum at lag 0 and decreases as it is shifted forward or backward.

## 2) Term 2



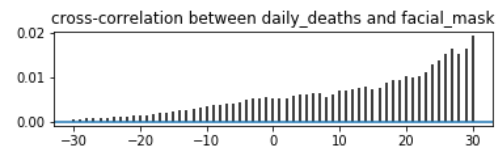
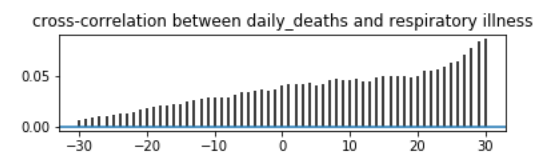
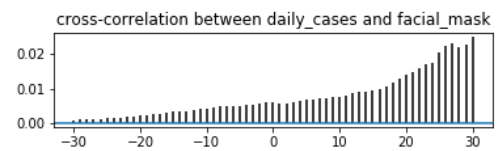
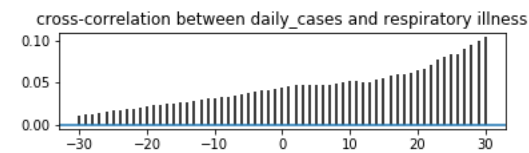
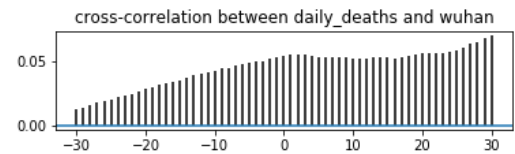
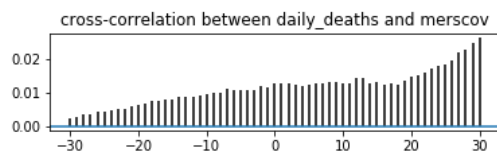
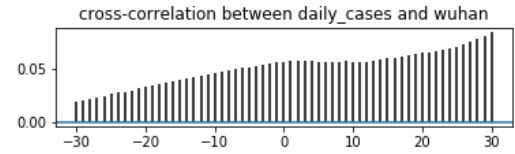
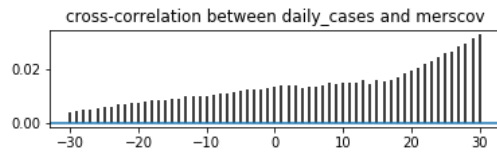
## 3) Term 3



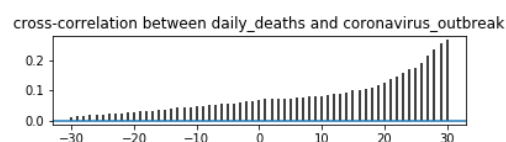
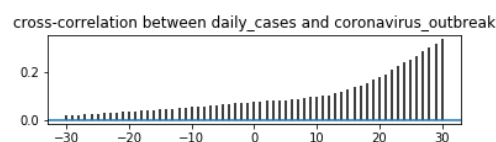
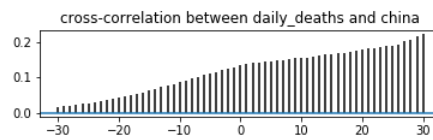
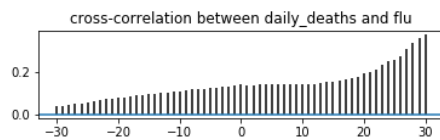
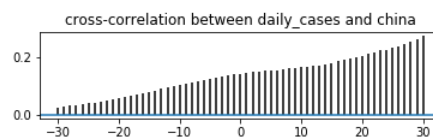
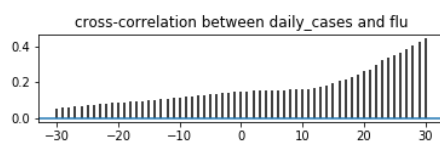
The above terms show higher correlation at higher lags, which means that when time series plots of the above are shifted right, they correlate well with “infection” and “deaths”.

Thus, to predict infection or death at time 't', counts of the above words at time 't-v', ('v' is a suitable lag ) can be better predictor.

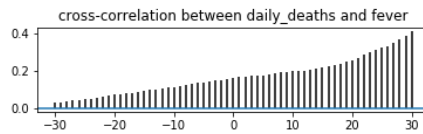
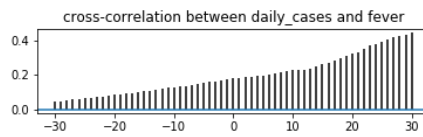
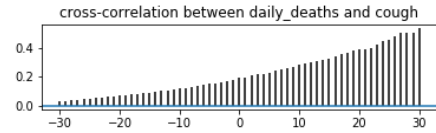
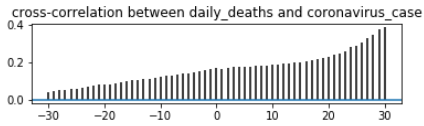
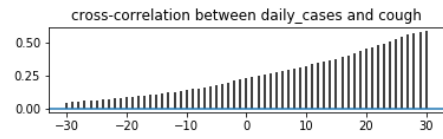
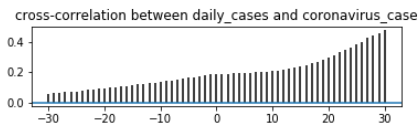
#### 4) Term 4



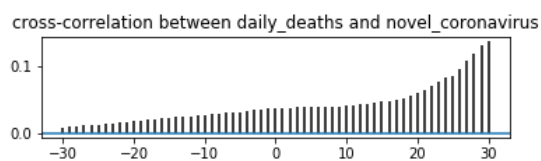
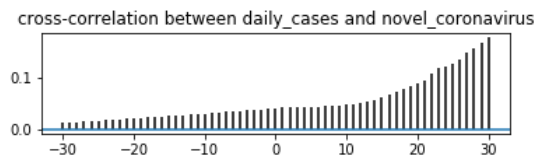
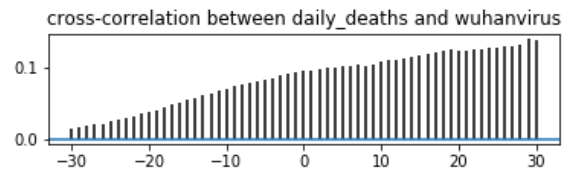
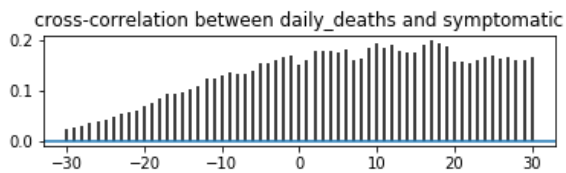
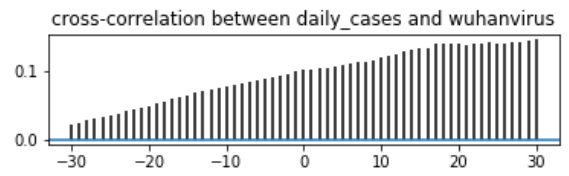
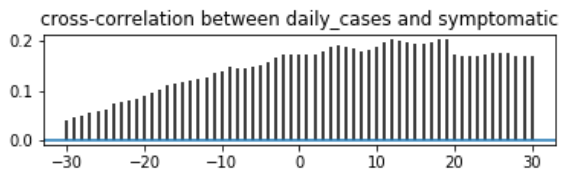
#### 5) Term 5



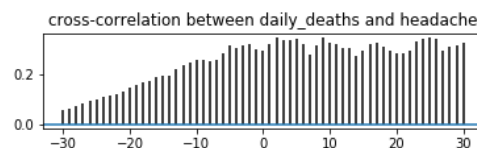
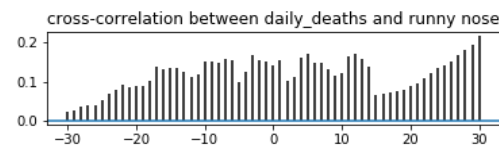
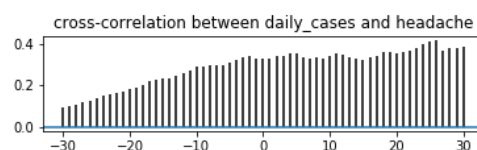
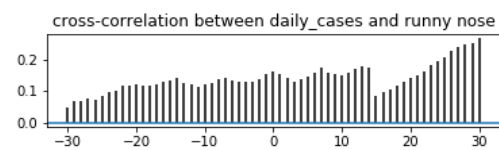
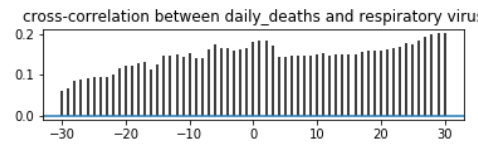
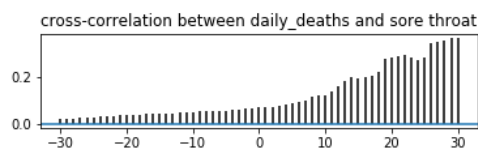
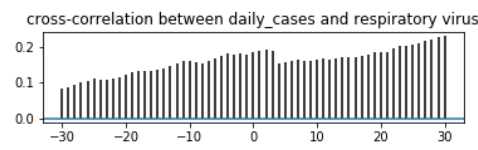
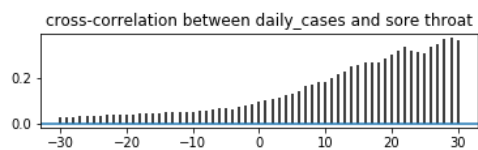
## Term 6



## 7) Term 7



8) Term 8



## Time series plots of clubbed terms

