

# Heuristic exploration of the relationship between FirstConfirmed.Per100K and viral load

Marlin Lee

1/28/2022

**This report looks is an update to the analysis shown on 1/14/2022. Most steps are the same with tweaking due to different data and outlier procedure**

at exploring the relationship between wastewater and FirstConfirmed.Per100K. There are four components to this analysis.

- 1) Removing putative outliers
- 2) Binning analysis
- 3) Smoothing signal
- 4) Statistical analysis

This report does not present any final answers but presents some very convincing heuristics.

“data Used from DSIWastewater package”

## Data: The first look

The two data sets used in this analysis are the Madison case data sourced from the Wisconsin DHS and wastewater concentration data produced by the Wisconsin State Laboratory of Hygiene. This wastewater data has entries every couple of days from 15 September 2020 to 19 April 2022.

```
## # A tibble: 1 x 3
##   site      'min(date)' 'max(date)'
##   <chr>    <date>        <date>
## 1 Madison 2020-09-15 2022-04-19
```

| date       | site    | FirstConfirmed.Per100K | pastwk.avg.casesperday.Per100K | sars_cov2_adj_load |
|------------|---------|------------------------|--------------------------------|--------------------|
| 2020-09-15 | Madison | 17.10526               | NA                             | 3.817117           |
| 2020-09-19 | Madison | 33.94737               | 37.40601                       | 1.578612           |
| 2020-09-22 | Madison | 17.10526               | 31.57895                       | 2.774853           |
| 2020-09-23 | Madison | 32.10526               | 27.96992                       | 2.459163           |
| 2020-09-24 | Madison | 23.68421               | 24.96241                       | 1.401820           |
| 2020-09-25 | Madison | 24.21053               | 23.00752                       | 2.345897           |

The case data has a strong weekend effect so for this section we look at a seven day smoothing of FirstConfirmed.Per100K. The simple display of the data shows the core components of this story. First, wastewater data is noisy. And that there is a clear relationship between the two signals.

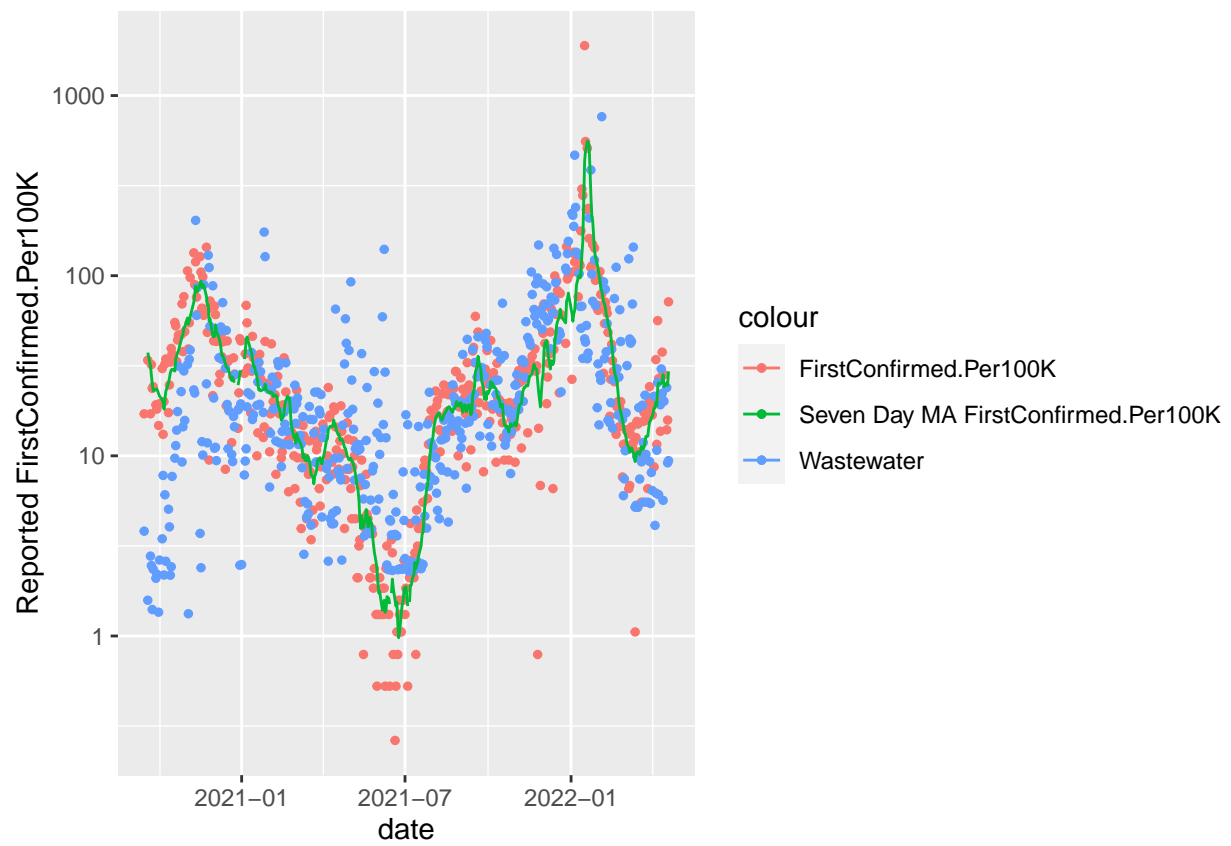


Figure 1: Wastewater concentration and daily Covid-19 case data for Madison. A seven day moving average of FirstConfirmed.Per100K is used to reduce a day of the week effect.

## Removing potential outliers

Looking at the wastewater measurements we observe there were some points many times larger than adjacent values hinting at them being outliers. We used the adjacent 10 values on each side and marked points 2.5 standard deviations away from the group mean as outliers.

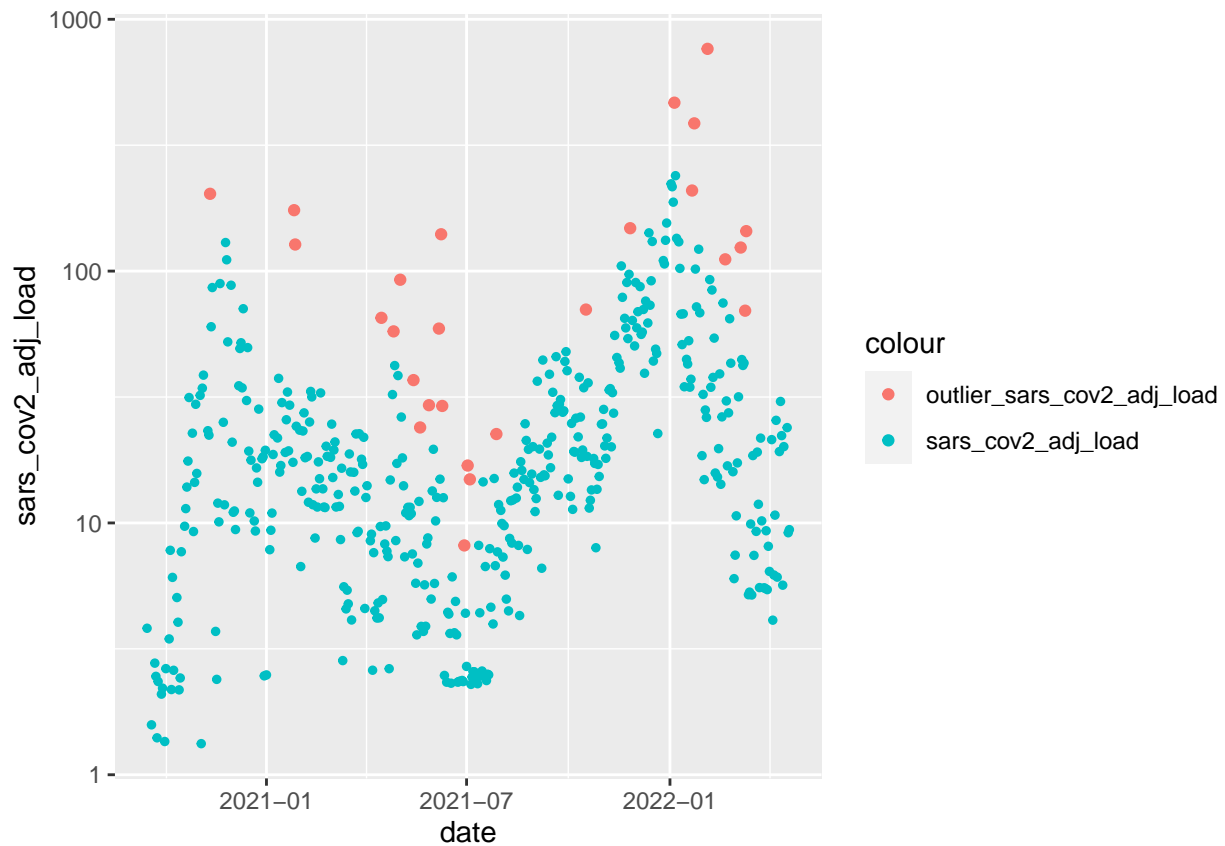


Figure 2: Wastewater concentration for Madison with potential outliers marked. Using a rolling symmetrical bin of 21 days as a sample we use 2.5 standard deviations of the bin as a metric to reject extreme points. This process is ran multiple times to get a robust process to select outliers.

## Data smoothing

The goal in this section is to smooth the data to get a similar effect without losing resolution.

### viral load smoothing

To get a good smoothing of the `sars_cov2_adj_load` measurement we employ loess smoothing. Loess smoothing takes a locally weighted sliding window using some number of points. we found the best smoothing when it uses data within approximately 0 weeks of both sides of the data. The displayed plot shows the visual power of this smoothing. We see in general that the smoothed N1 trails SLD. However loess is symmetric meaning that it can not be used in predictive modeling due to it using points from the future to smooth points.



Figure 3: Loess smoothed N1 and SLD FirstConfirmed.Per100K for Madison data. Using a Locally Weighted Scatterplot Smoothing process along with the previous figure SLD FirstConfirmed.Per100K we get the most sophisticated relationship between the two signals discussed in this document.

## Towards a formal analysis

Cross correlation and Granger Causality are key components to formalize this analysis. Cross correlation looks at the correlation at a range of time shifts and Granger analysis performs a test for predictive power.

|  | Max<br>Cross<br>Correla-<br>tion | Lag of<br>largest Cross<br>correlation | P-value Wastewater<br>predicts FirstCon-<br>firmed.Per100K | P-value FirstCon-<br>firmed.Per100K<br>predicts wastewater |
|--|----------------------------------|--|--|--|
| Section 1:<br>FirstConfirmed.Per100K vs<br>sars_cov2_adj_load                                  | 0.4761                           | 18                                     | 0.0433   | 0.0000   |
| Section 1: 7 Day MA<br>FirstConfirmed.Per100K vs<br>sars_cov2_adj_load                         | 0.5436                           | 13                                     | 0.0167   | 0.0000   |
| Section 2:<br>FirstConfirmed.Per100K vs<br>sars_cov2_adj_load                                  | 0.5899                           | 15                                     | 0.3797   | 0.0001   |
| Section 4.3: 7 Day MA<br>FirstConfirmed.Per100K vs Loess<br>smoothing of<br>sars_cov2_adj_load | 0.7279                           | 1                                      | 0.0314   | 0.0767   |