

# Exploration of COVID-19 tracking data from multiple resources

Wei Sun

2020-04-25

## Contents

<b>Introduction</b>	<b>1</b>
<b>JHU</b>	<b>2</b>
time series data . . . . .	2
daily reports data . . . . .	6
<b>NY Times</b>	<b>7</b>
state level data . . . . .	7
county level data . . . . .	14
<b>COVID Trackng</b>	<b>21</b>
<b>Session information</b>	<b>22</b>

## Introduction

Coronavirus disease 2019 (COVID-19) is an infectious disease caused by a new type of coronavirus: severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2). The outbreak first started in Wuhan, China in December 2019. The first kown case of COVID-19 in the U.S. was confirmed on January 20, 2020, in a 35-year-old man who teturned to Washington State on January 15 after traveling to Wuhan. Starting around the end of Feburary, evidence emerge for community spread in the US.

We, as all of us, are indebted to the heros who fight COVID-19 across the whole world in different ways. For this data exploration, I am grateful to many data science groups who have collected detailed COVID-19 outbreak data, including the number of tests, confirmed cases, and deaths, across countries/regions, states/provnices (administrative division level 1, or admin1), and counties (admin2). Specifically, I used the data from these three resources:

- JHU (<https://coronavirus.jhu.edu/>)
  - The Center for Systems Science and Engineering (CSSE) at John Hopkins University.
  - World-wide counts of coronavirus cases, deaths, and recovered ones.
  - <https://github.com/CSSEGISandData/COVID-19>
- NY Times (<https://www.nytimes.com/interactive/2020/us/coronavirus-us-cases.html>)
  - The New York Times
  - “cumulative counts of coronavirus cases in the United States, at the state and county level, over time”
  - <https://github.com/nytimes/covid-19-data>

- COVID Tracking (<https://covidtracking.com/>)
  - COVID Tracking Project
  - “collects information from 50 US states, the District of Columbia, and 5 other US territories to provide the most comprehensive testing data”
  - <https://github.com/COVID19Tracking/covid-tracking-data>

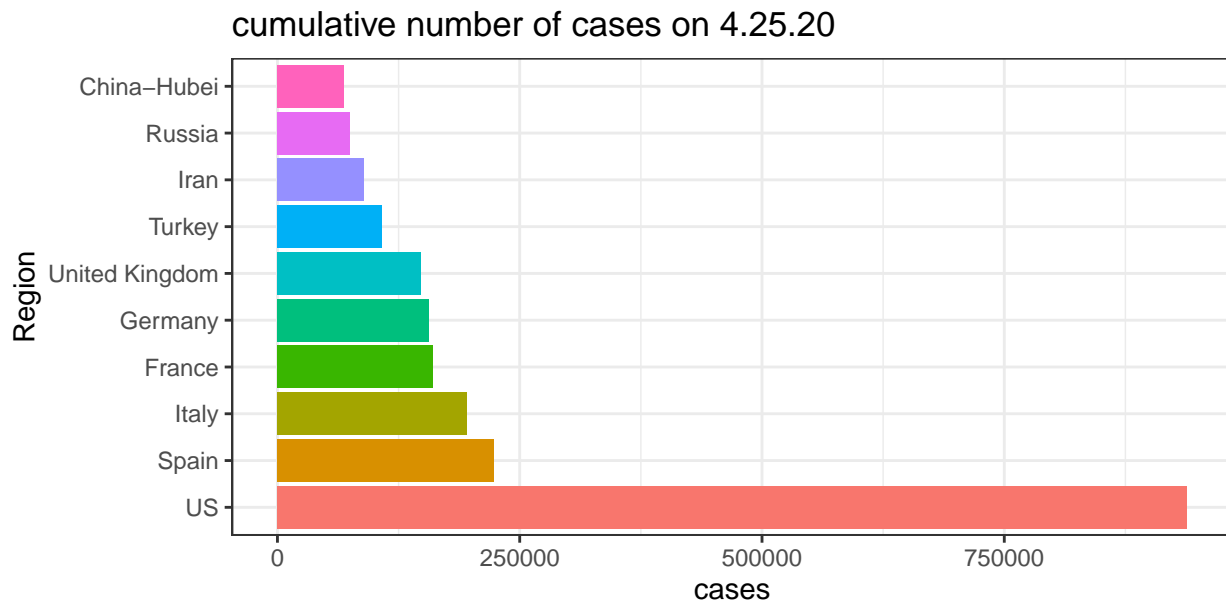
## JHU

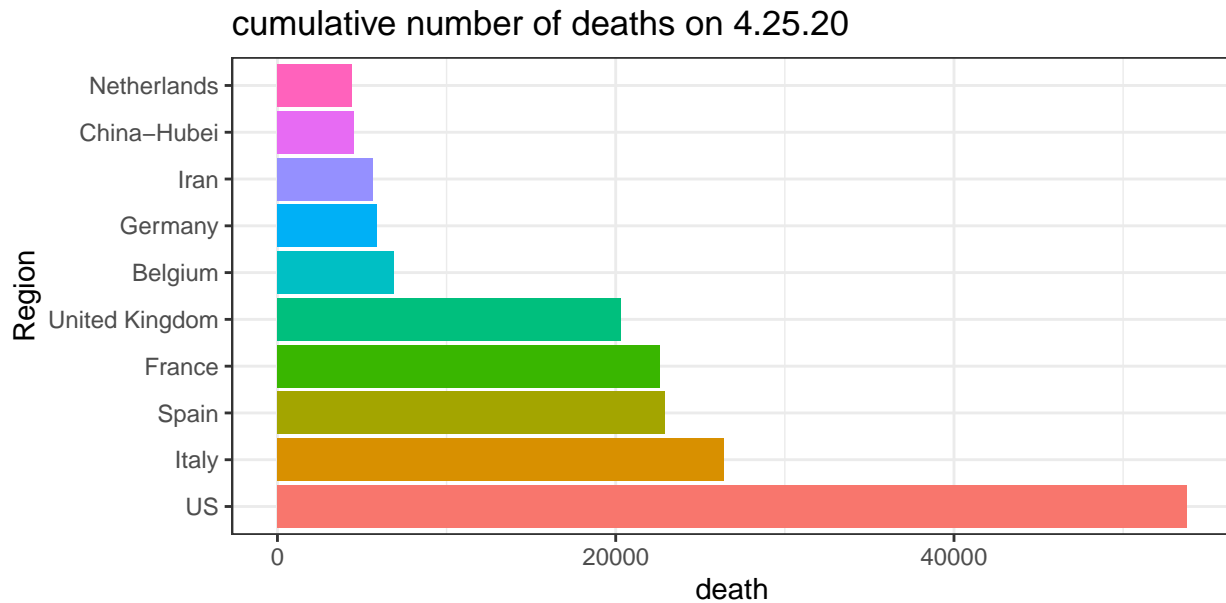
Assume you have cloned the JHU Github repository on your local machine at “../COVID-19”.

### time series data

The time series provide counts (e.g., confirmed cases, deaths) starting from Jan 22nd, 2020 for 253 locations. Currently there is no data of individual US state in these time series data files.

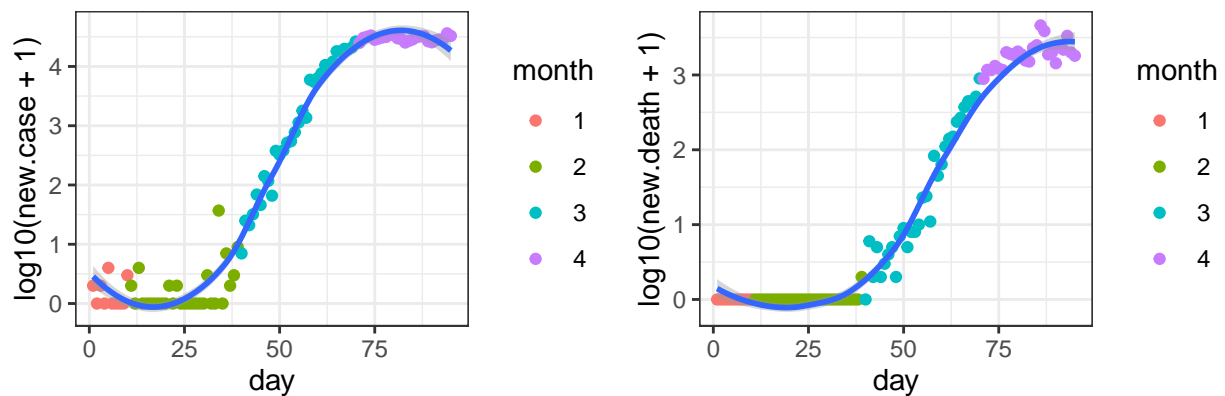
Here is the list of 10 records with the largest number of cases or deaths on the most recent date.





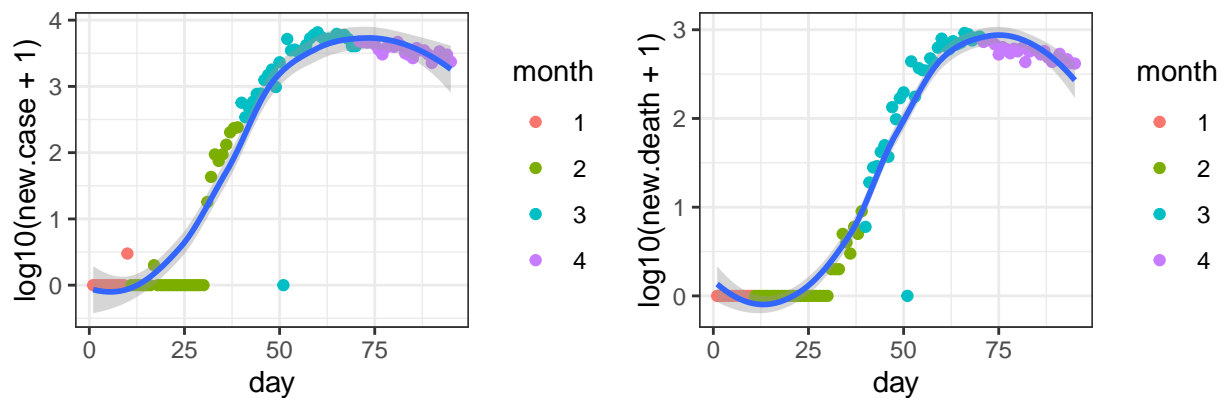
Next, I check for each country/region, what is the number of new cases/deaths? This data is important to understand what is the trend under different situations, e.g., population density, social distance policies etc. Here I checked the top 10 countries/regions with the highest number of deaths.

### US



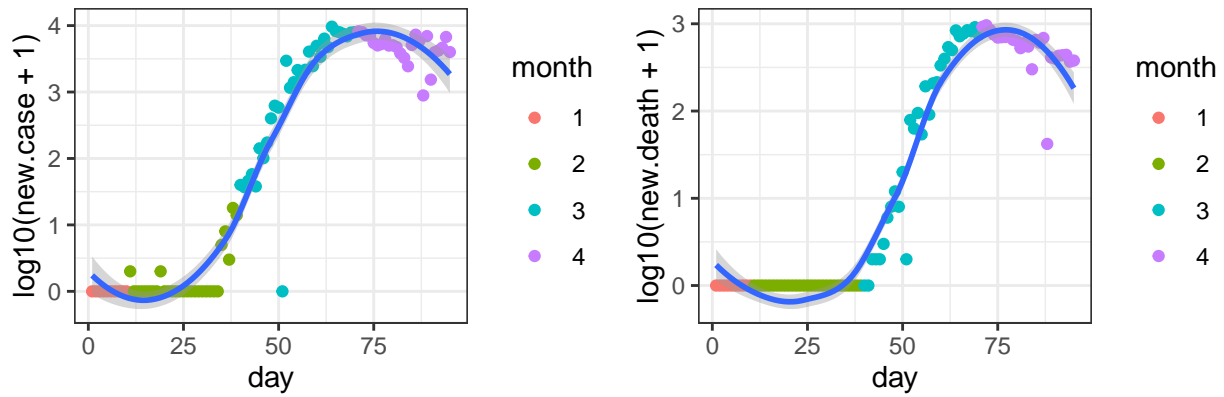
data source: <https://github.com/CSSEGISandData/COVID-19>, day 1 is 1/22/2020

### Italy



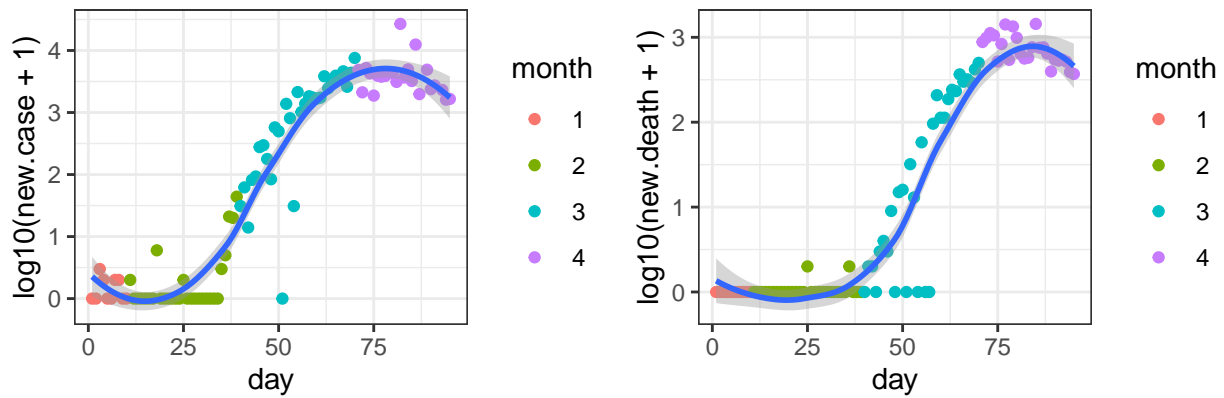
data source: <https://github.com/CSSEGISandData/COVID-19>, day 1 is 1/22/2020

### Spain



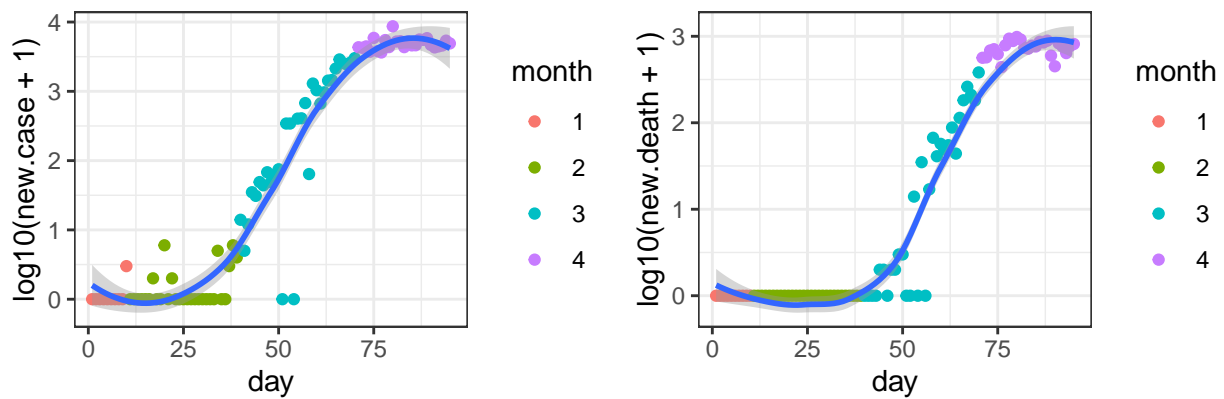
data source: <https://github.com/CSSEGISandData/COVID-19>, day 1 is 1/22/2020

### France



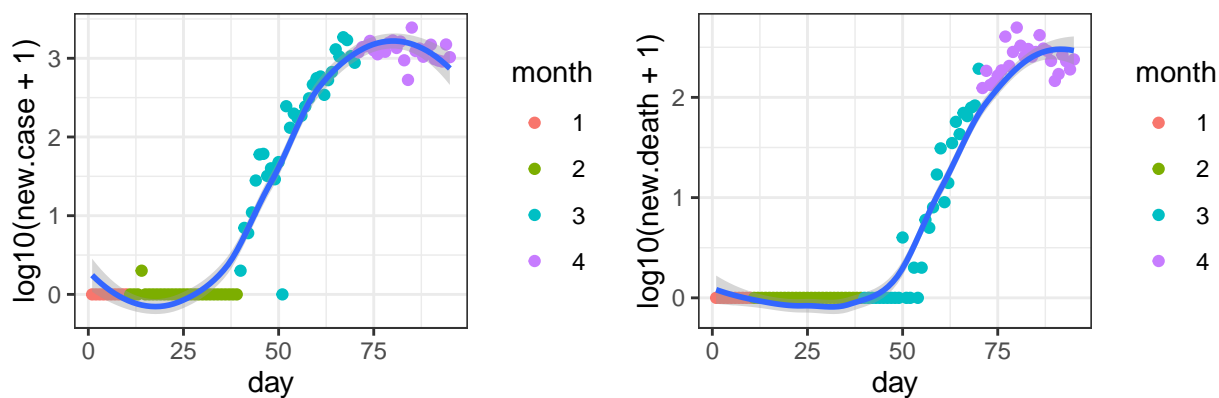
data source: <https://github.com/CSSEGISandData/COVID-19>, day 1 is 1/22/2020

### United Kingdom



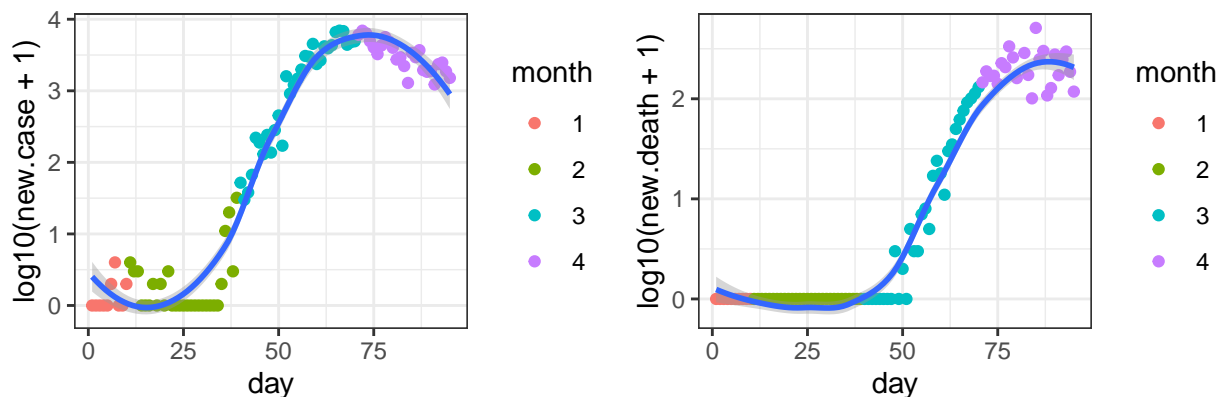
data source: <https://github.com/CSSEGISandData/COVID-19>, day 1 is 1/22/2020

## Belgium



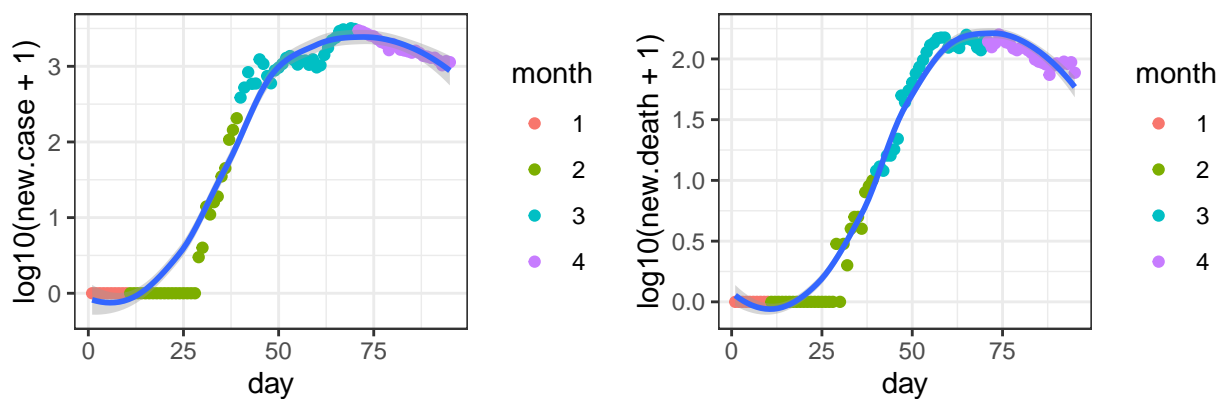
data source: <https://github.com/CSSEGISandData/COVID-19>, day 1 is 1/22/2020

## Germany



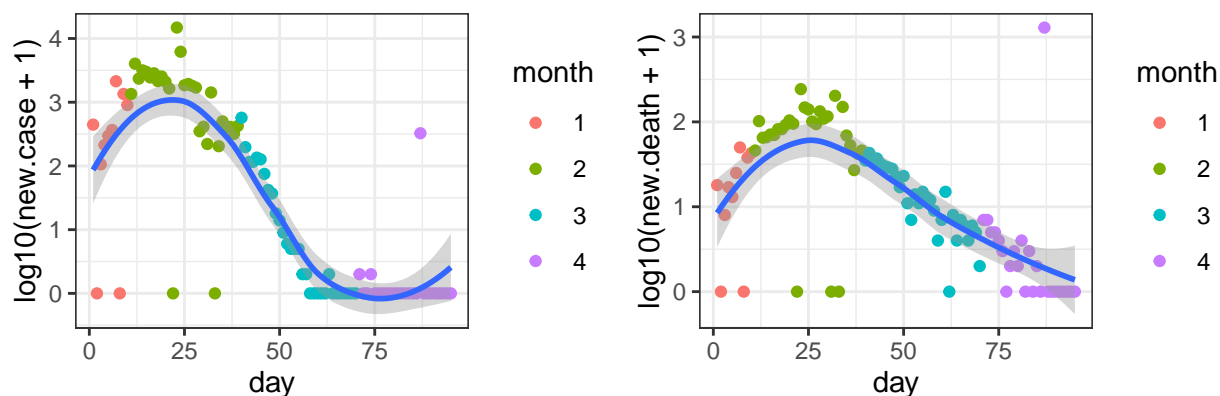
data source: <https://github.com/CSSEGISandData/COVID-19>, day 1 is 1/22/2020

## Iran



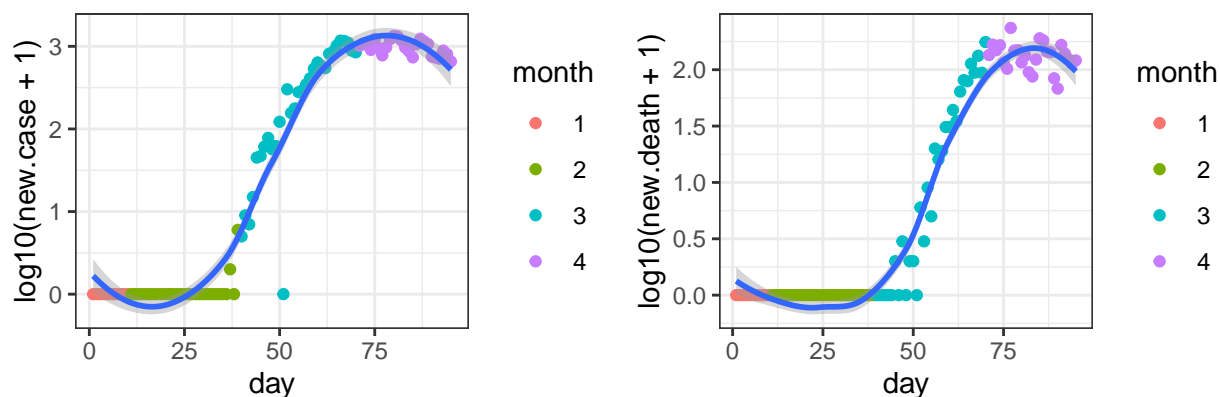
data source: <https://github.com/CSSEGISandData/COVID-19>, day 1 is 1/22/2020

## China-Hubei



data source: <https://github.com/CSSEGISandData/COVID-19>, day 1 is 1/22/2020

## Netherlands

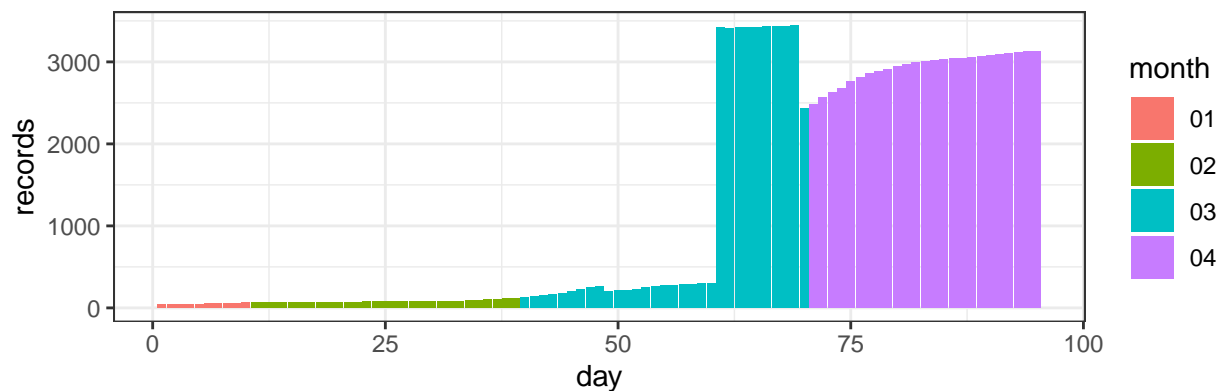


data source: <https://github.com/CSSEGISandData/COVID-19>, day 1 is 1/22/2020

## daily reports data

The raw data from Hopkins are in the format of daily reports with one file per day. More recent files (since March 22nd) include information from individual states of US or individual counties, as shown in the following figure. So I turn to NY Times data for informatoin of individual states or counties.

### number of records in Hopkins daily reports



data source: <https://github.com/CSSEGISandData/COVID-19>, day 1 is 1/22/2020

## NY Times

The data from NY Times are saved in two text files, one for state level information and the other one for county level information.

The current date is

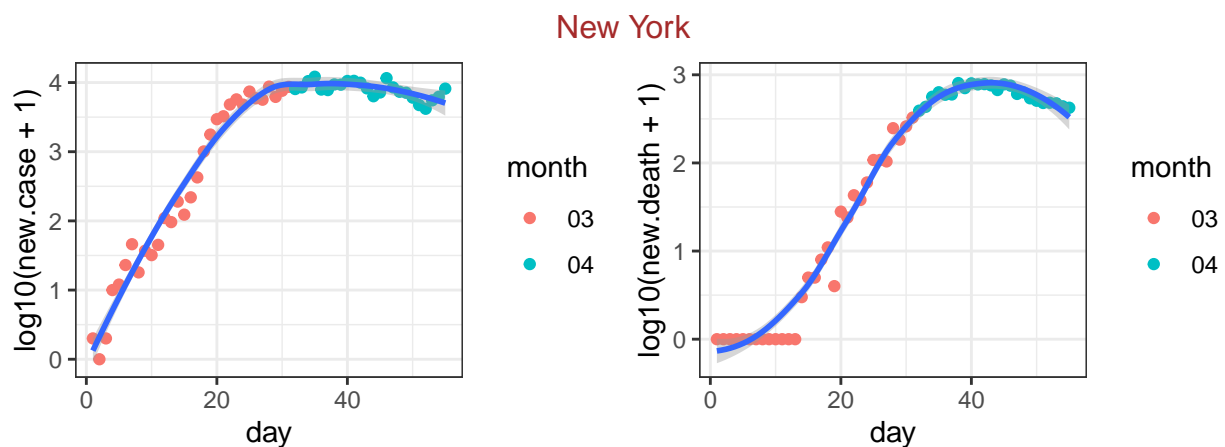
```
## [1] "2020-04-24"
```

### state level data

First check the 20 states with the largest number of deaths.

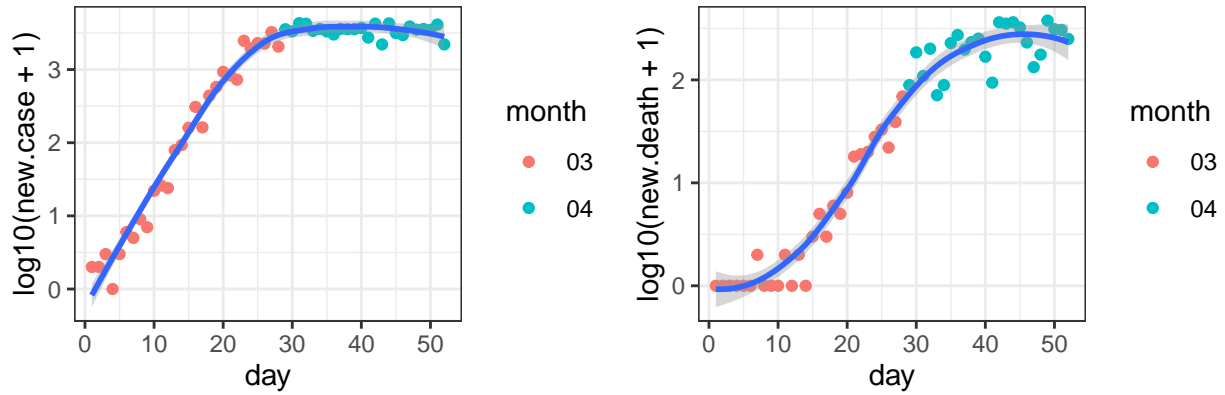
##	date	state	fips	cases	deaths	
##	2908	2020-04-24	New York	36	271621	16162
##	2906	2020-04-24	New Jersey	34	102196	5617
##	2898	2020-04-24	Michigan	26	36627	3084
##	2897	2020-04-24	Massachusetts	25	50969	2556
##	2889	2020-04-24	Illinois	17	39658	1804
##	2915	2020-04-24	Pennsylvania	42	40298	1786
##	2881	2020-04-24	Connecticut	9	23921	1764
##	2879	2020-04-24	California	6	41368	1619
##	2894	2020-04-24	Louisiana	22	26140	1601
##	2884	2020-04-24	Florida	12	30525	1045
##	2885	2020-04-24	Georgia	13	21575	889
##	2890	2020-04-24	Indiana	18	13680	741
##	2926	2020-04-24	Washington	53	13120	731
##	2896	2020-04-24	Maryland	24	16618	723
##	2912	2020-04-24	Ohio	39	15169	690
##	2880	2020-04-24	Colorado	8	12255	672
##	2921	2020-04-24	Texas	48	23650	625
##	2925	2020-04-24	Virginia	51	11596	413
##	2909	2020-04-24	North Carolina	37	8052	270
##	2877	2020-04-24	Arizona	4	6045	268

For these 20 states, I check the number of new cases and the number of new deaths. Part of the reason for such checking is to identify whether there is any similarity on such patterns. For example, could you use the pattern seen from Italy to predict what happen in an individual state, and what are the similarities and differences across states.



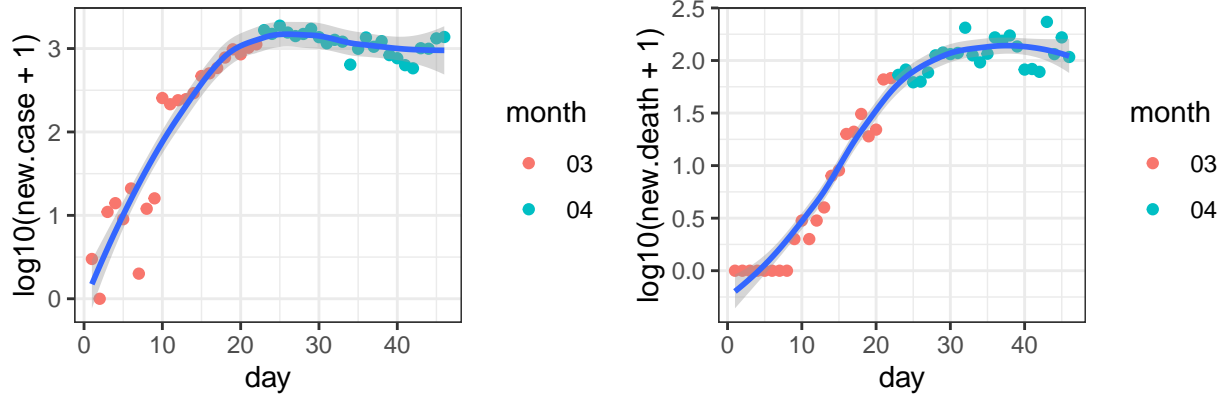
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-01

### New Jersey



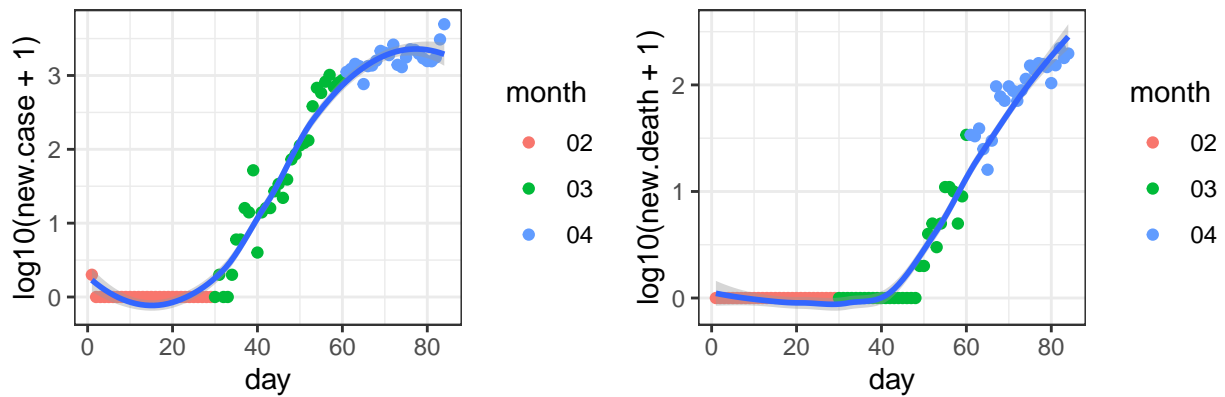
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-04

### Michigan



data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-10

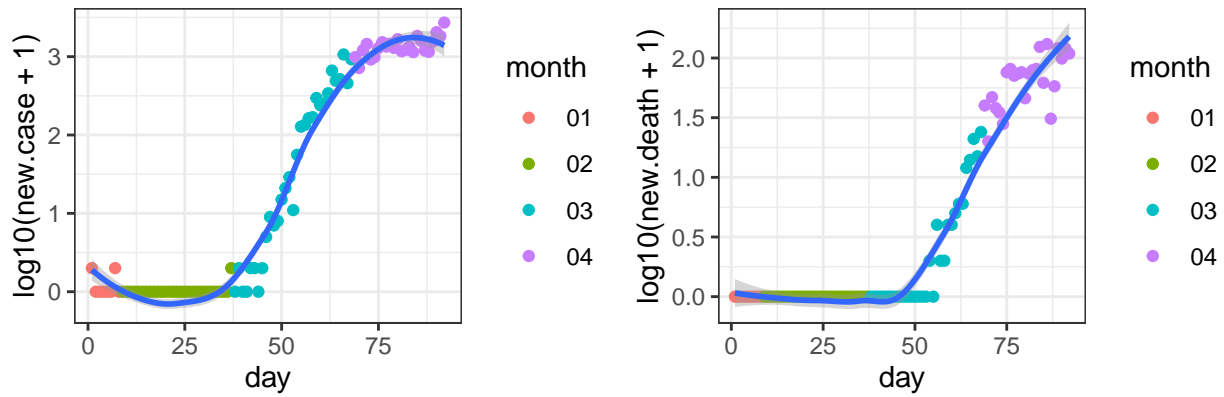
### Massachusetts



data source: <https://github.com/nytimes/covid-19-data>, day 1 is 02-01

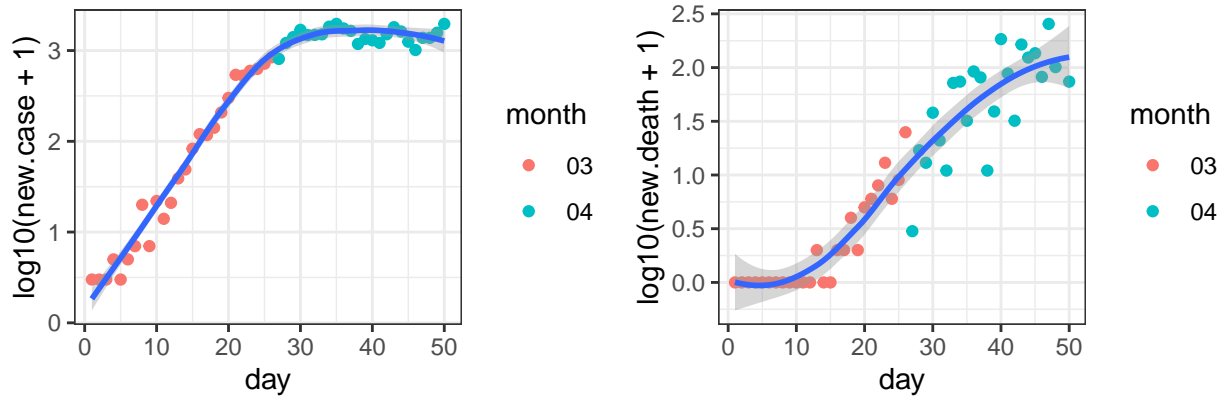


### Illinois



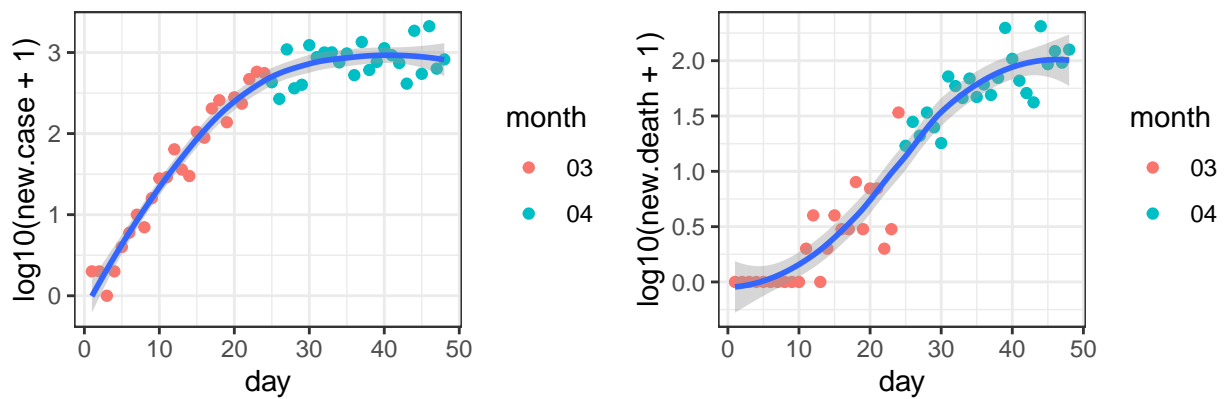
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 01-24

### Pennsylvania



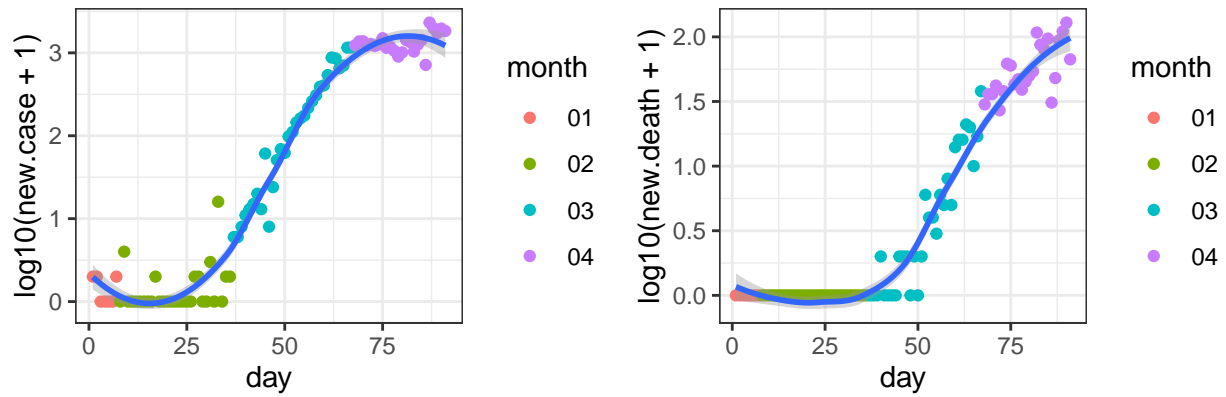
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-06

### Connecticut



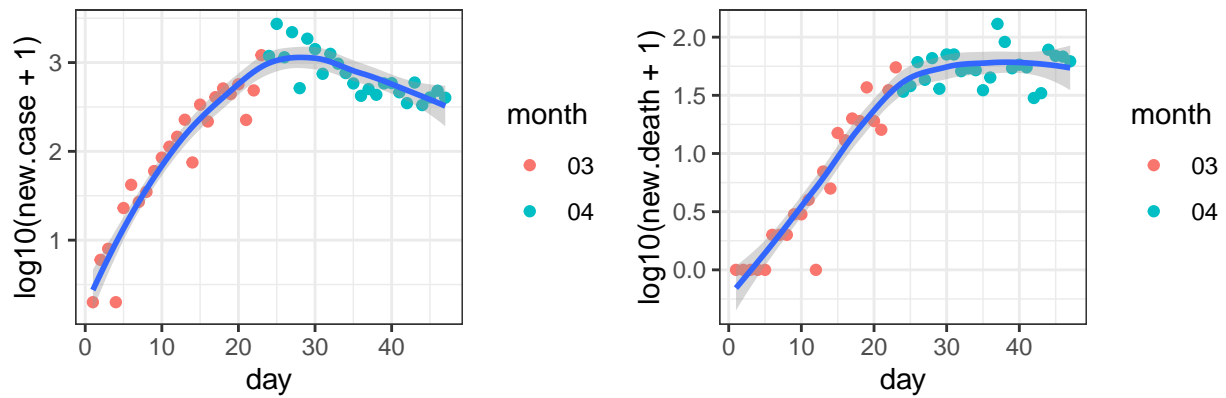
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-08

### California



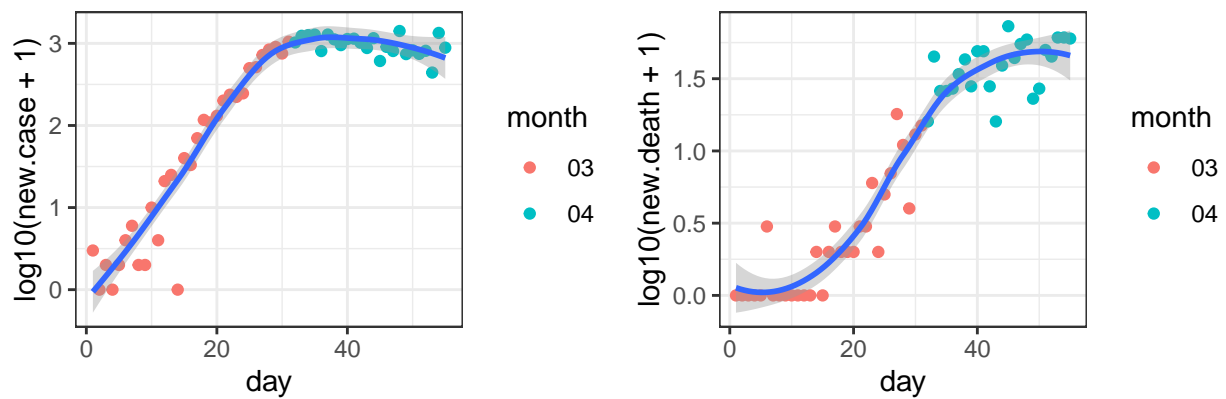
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 01-25

### Louisiana

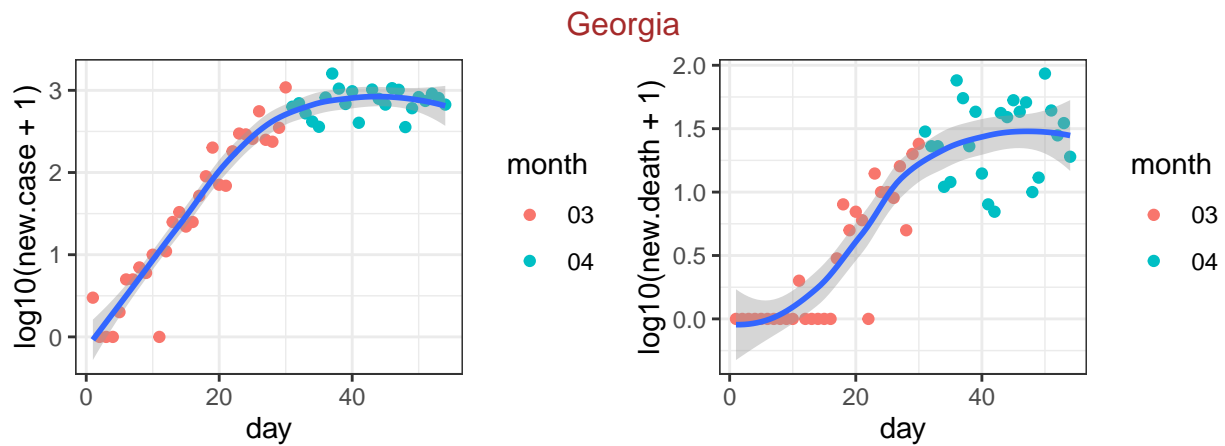


data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-09

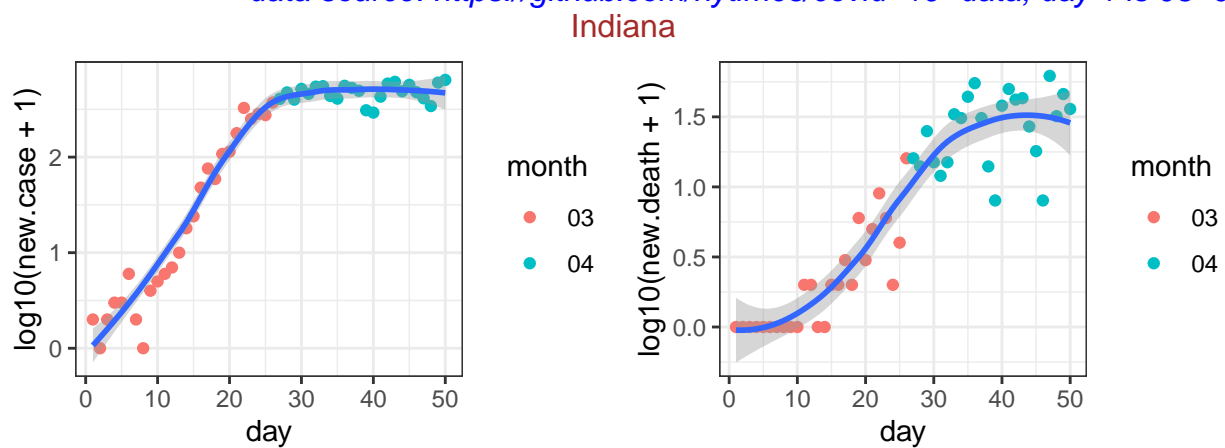
### Florida



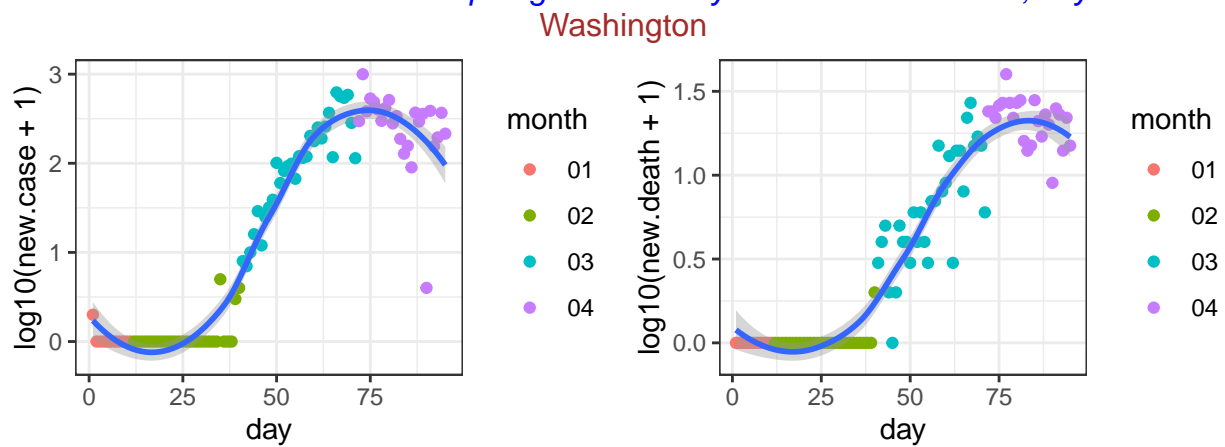
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-01



data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-02

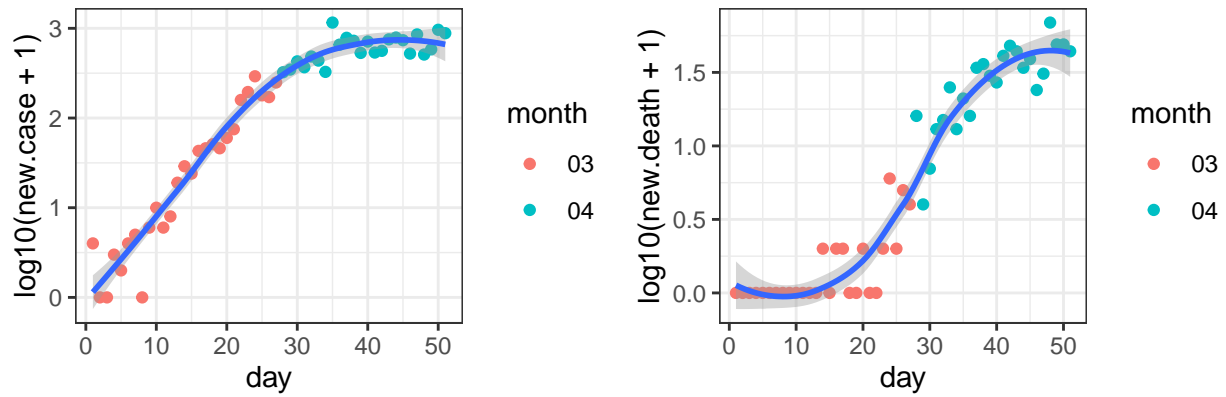


data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-06



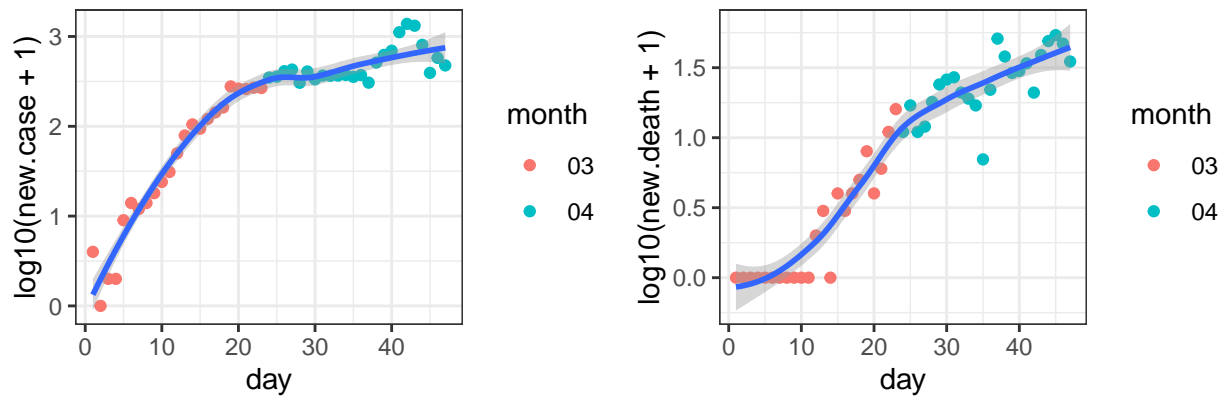
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 01-21

### Maryland



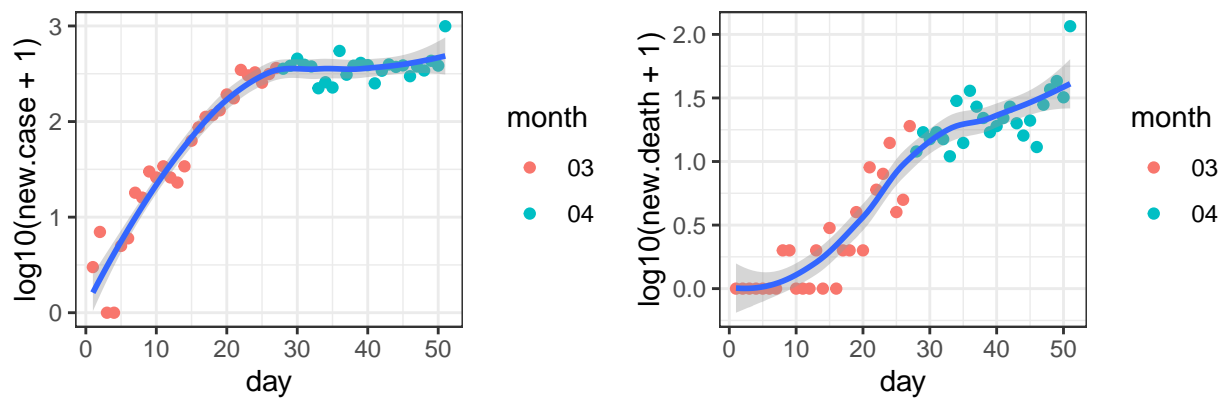
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-05

### Ohio



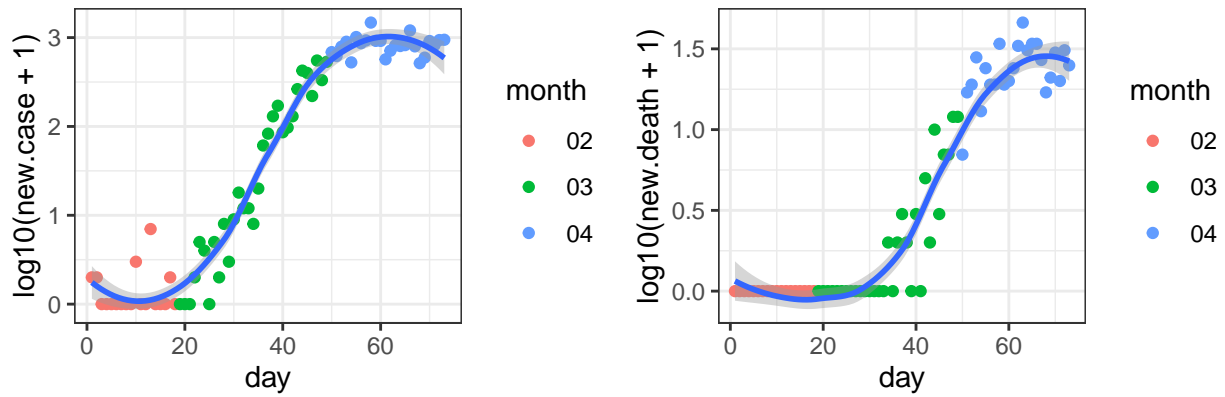
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-09

### Colorado



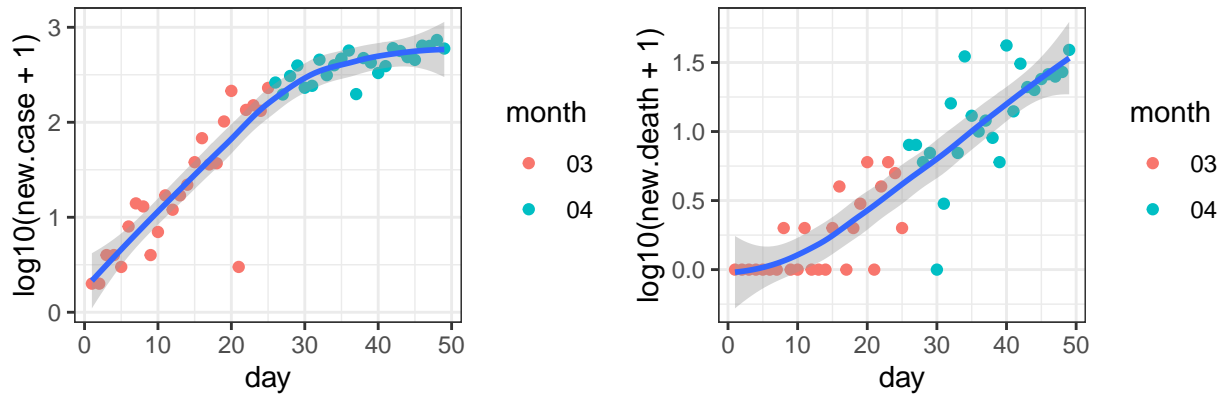
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-05

### Texas



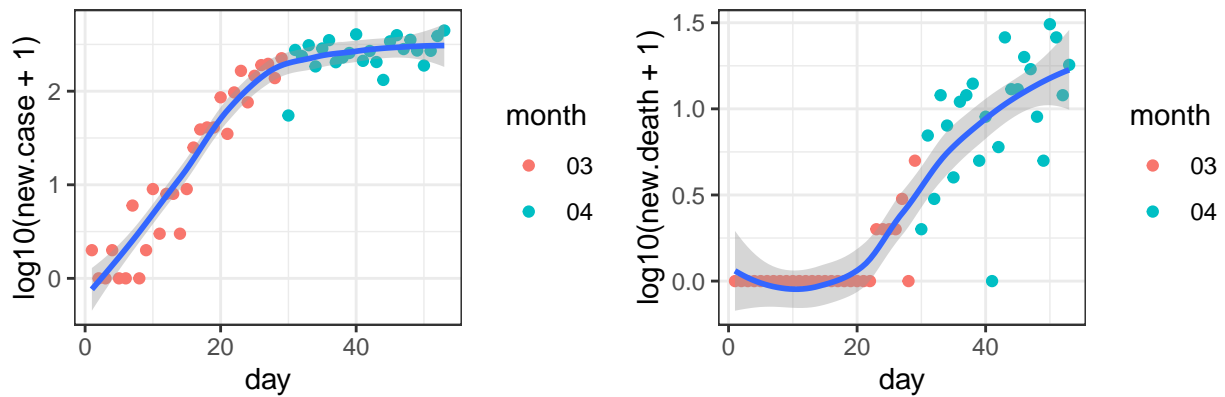
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 02-12

### Virginia

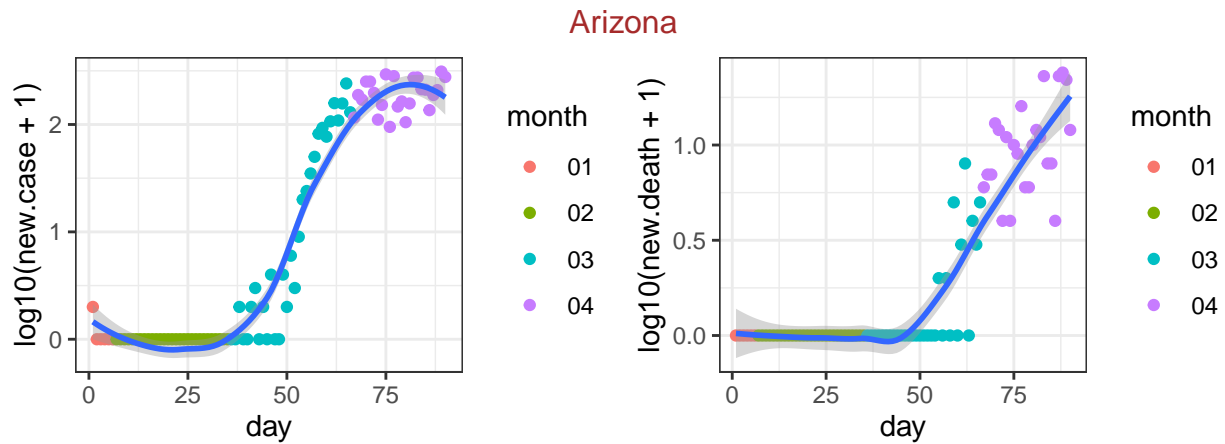


data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-07

### North Carolina

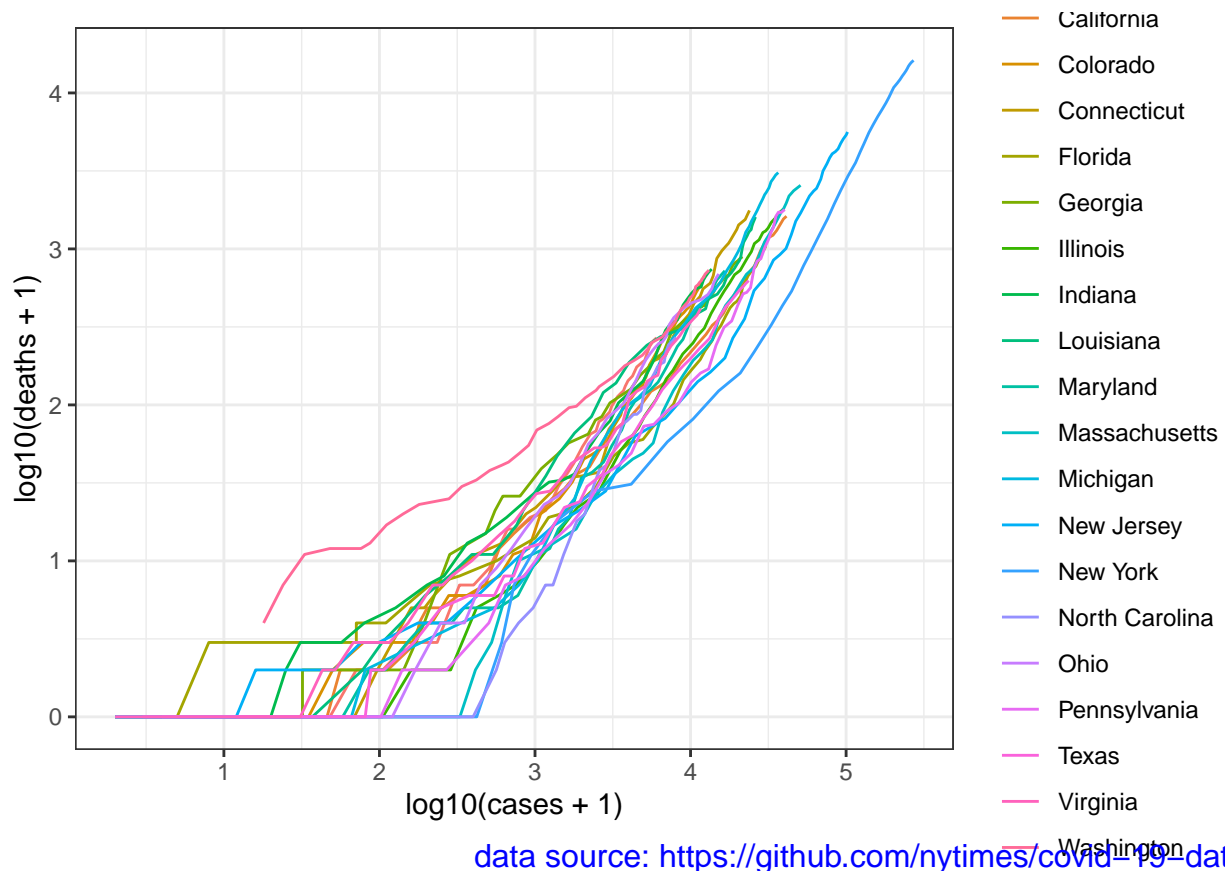


data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-03



data source: <https://github.com/nytimes/covid-19-data>, day 1 is 01-26

Next I check the relation between the **cumulative** number of cases and deaths for these 10 states, starting on March



data source: <https://github.com/nytimes/covid-19-data>

## county level data

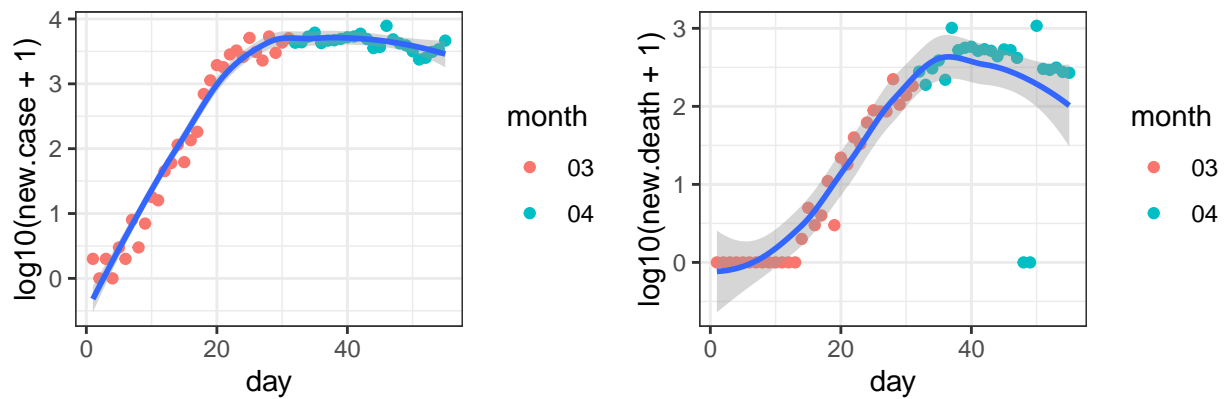
First check the 20 counties with the largest number of deaths.

##	date	county	state	fips	cases	deaths
## 85808	2020-04-24	New York City	New York	NA	150484	11157
## 85807	2020-04-24	Nassau	New York	36059	32765	1867
## 85359	2020-04-24	Wayne	Michigan	26163	15407	1443

##	84712	2020-04-24	Cook	Illinois	17031	27616	1220
##	85827	2020-04-24	Suffolk	New York	36103	30606	1035
##	85835	2020-04-24	Westchester	New York	36119	26632	989
##	85736	2020-04-24	Essex	New Jersey	34013	12110	975
##	85731	2020-04-24	Bergen	New Jersey	34003	14363	934
##	84329	2020-04-24	Los Angeles	California	6037	18545	850
##	84422	2020-04-24	Fairfield	Connecticut	9001	10227	662
##	85738	2020-04-24	Hudson	New Jersey	34017	13011	640
##	85274	2020-04-24	Middlesex	Massachusetts	25017	11681	585
##	85340	2020-04-24	Oakland	Michigan	26125	6804	585
##	85749	2020-04-24	Union	New Jersey	34039	11208	542
##	84423	2020-04-24	Hartford	Connecticut	9003	4570	511
##	85327	2020-04-24	Macomb	Michigan	26099	5022	504
##	86203	2020-04-24	Philadelphia	Pennsylvania	42101	11877	449
##	85741	2020-04-24	Middlesex	New Jersey	34023	9789	413
##	84426	2020-04-24	New Haven	Connecticut	9009	6286	396
##	86795	2020-04-24	King	Washington	53033	5691	393

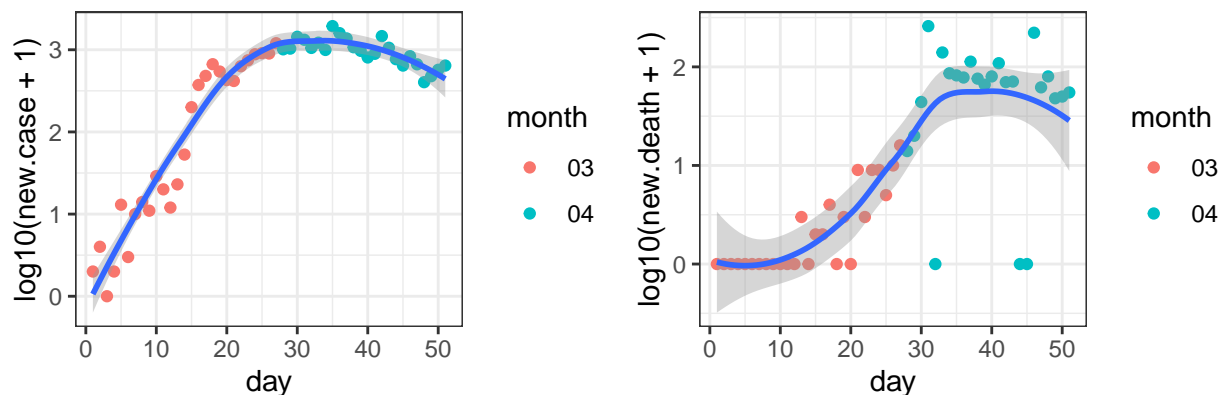
For these 20 counties, I check the number of new cases and the number of new deaths.

### New York City\_New York



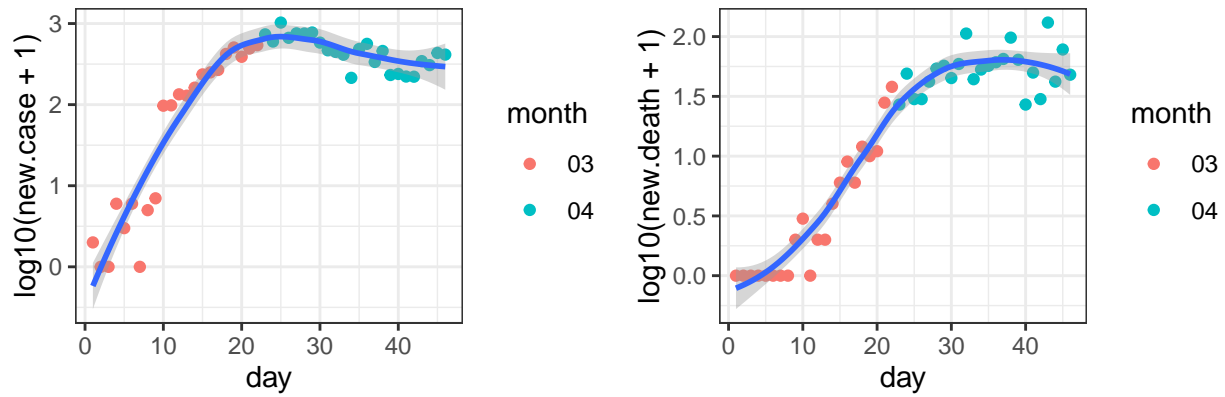
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-01

### Nassau\_New York



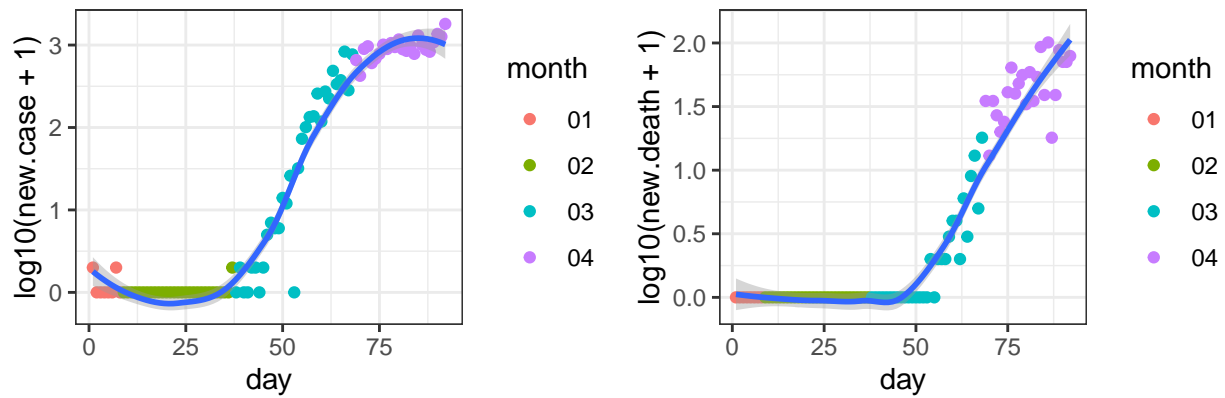
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-05

### Wayne\_Michigan



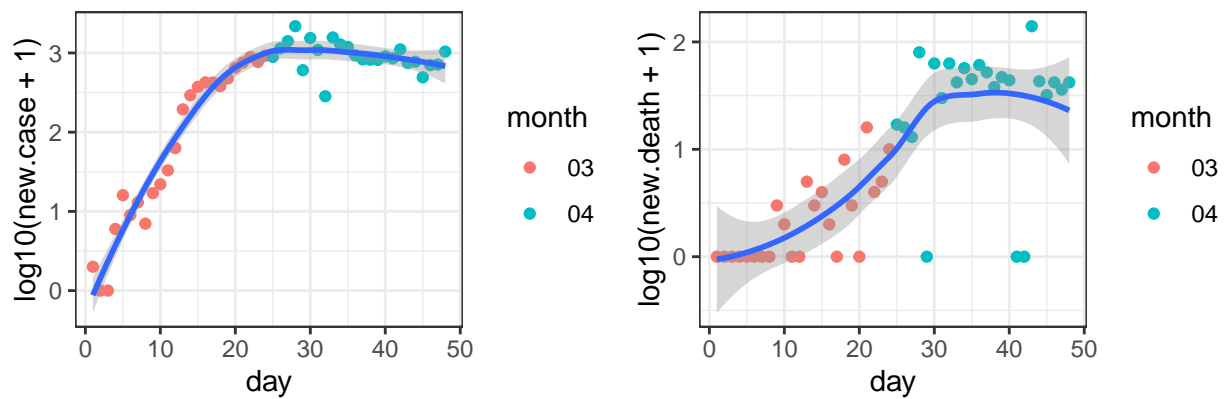
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-10

### Cook\_Illinois



data source: <https://github.com/nytimes/covid-19-data>, day 1 is 01-24

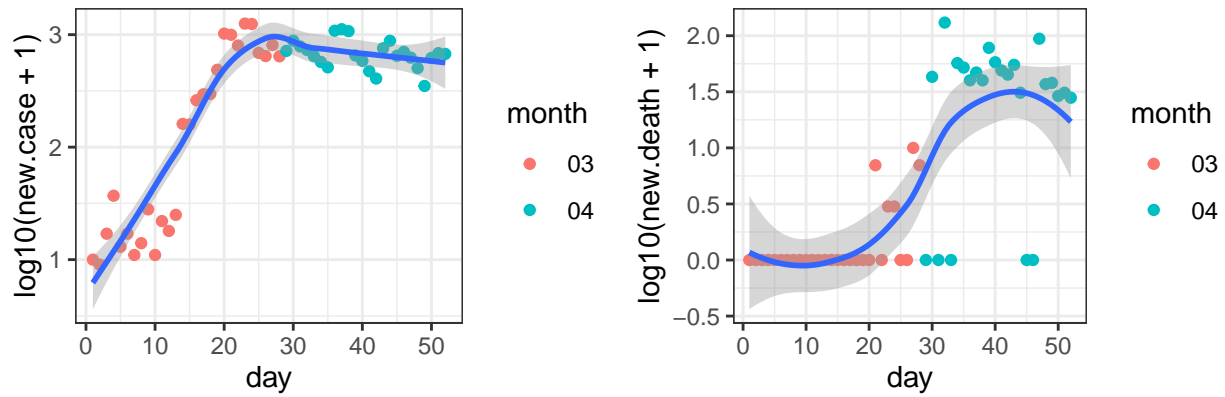
### Suffolk\_New York



data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-08

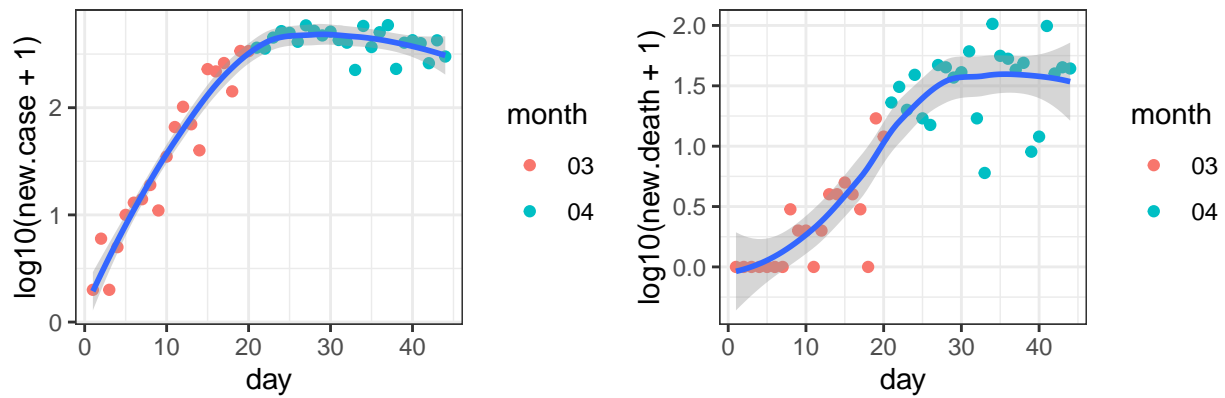


### Westchester\_New York



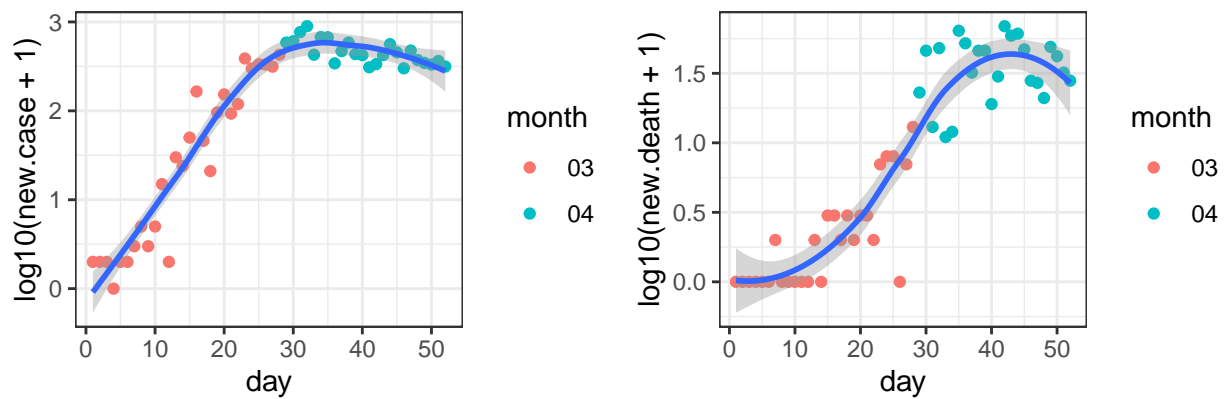
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-04

### Essex\_New Jersey



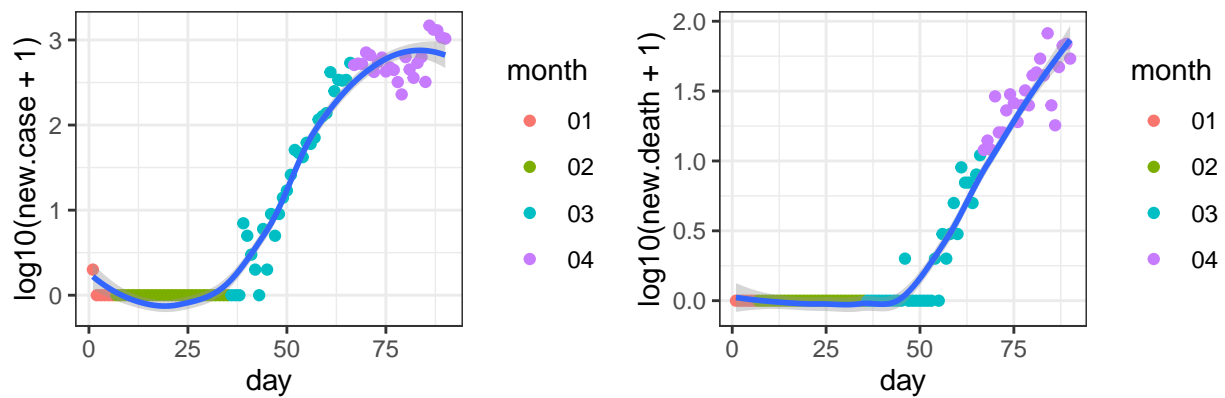
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-12

### Bergen\_New Jersey



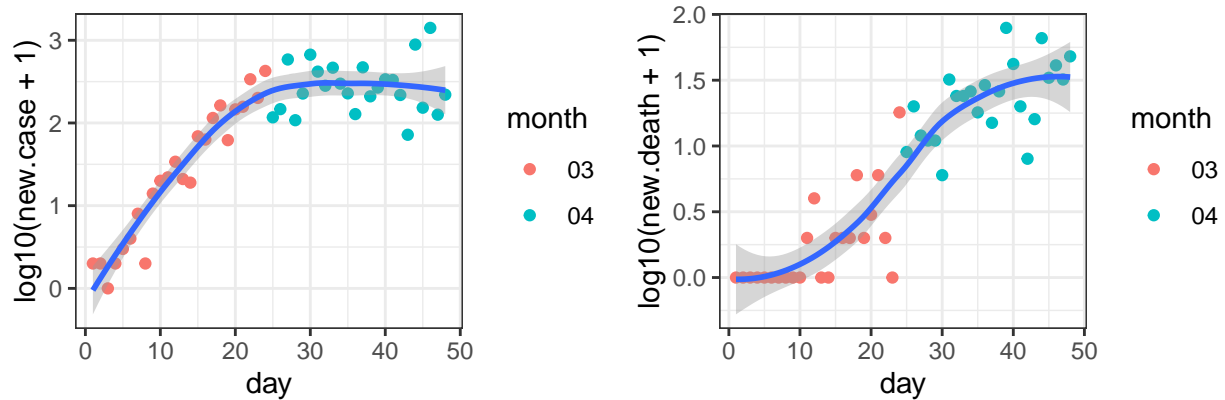
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-04

### Los Angeles\_California



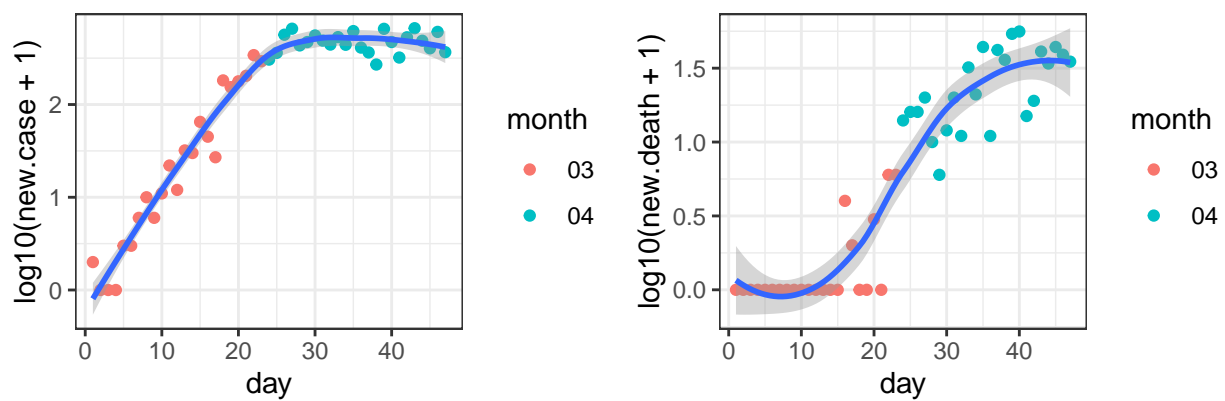
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 01-26

### Fairfield\_Connecticut



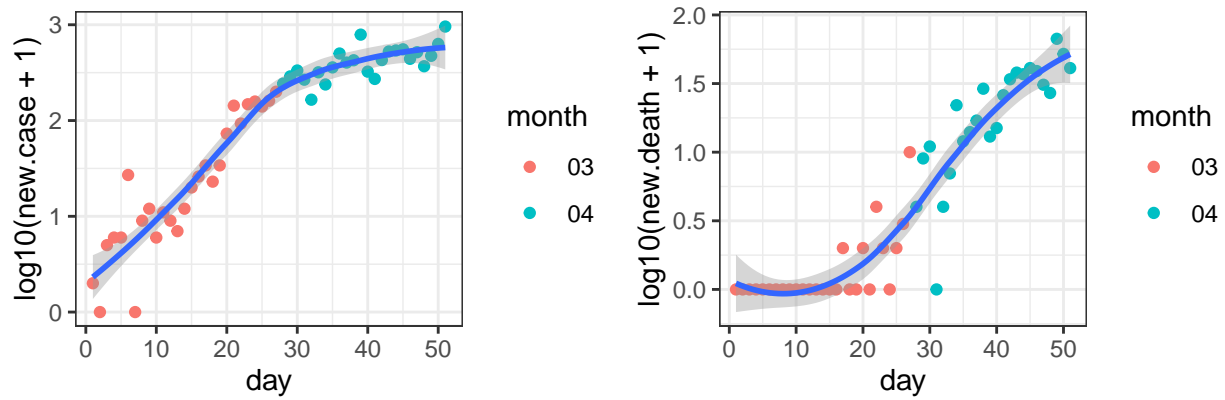
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-08

### Hudson\_New Jersey



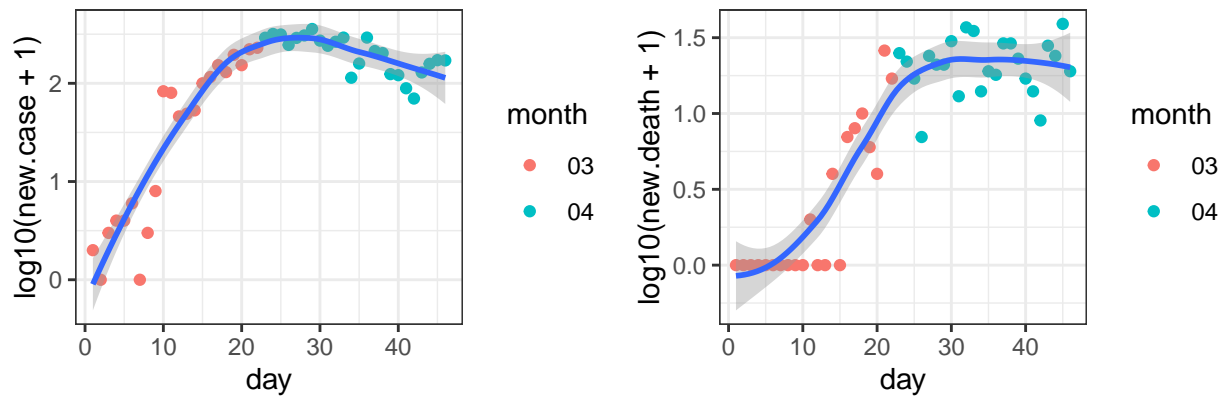
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-09

### Middlesex\_Massachusetts



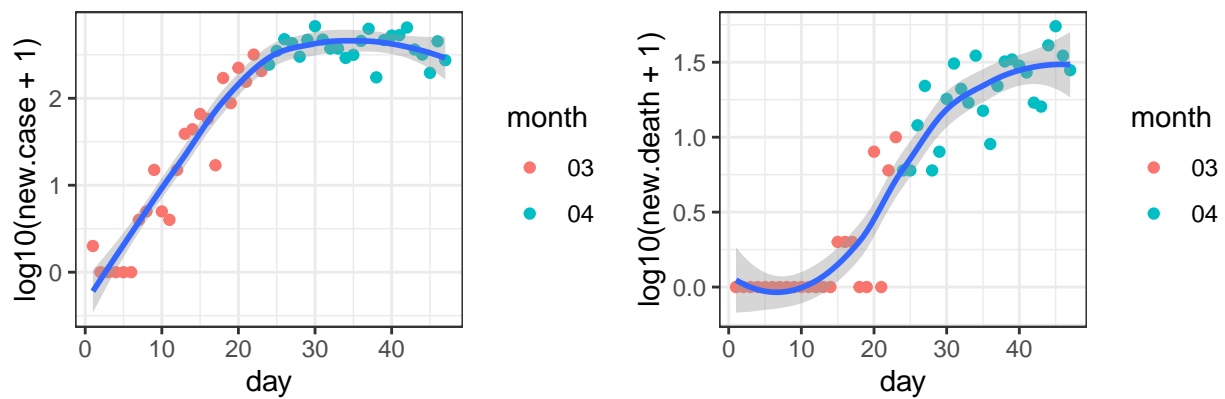
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-05

### Oakland\_Michigan



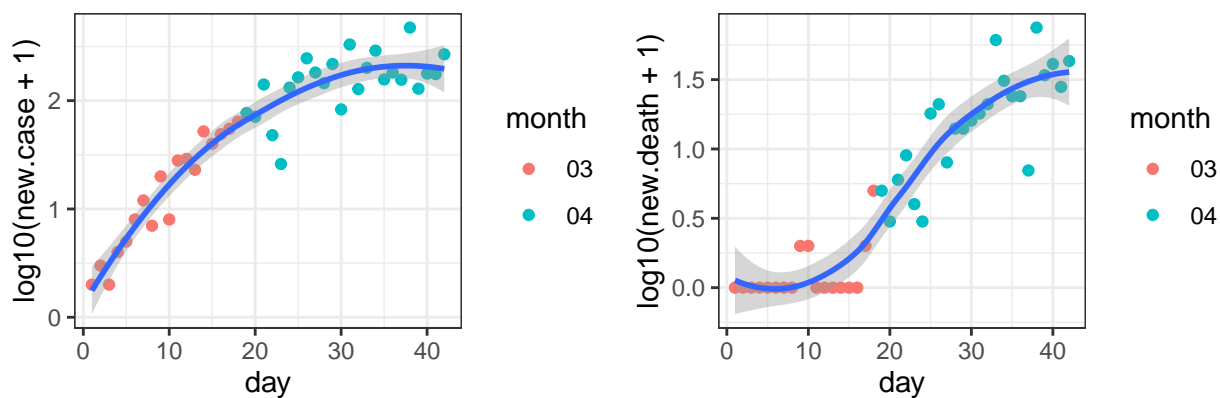
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-10

### Union\_New Jersey



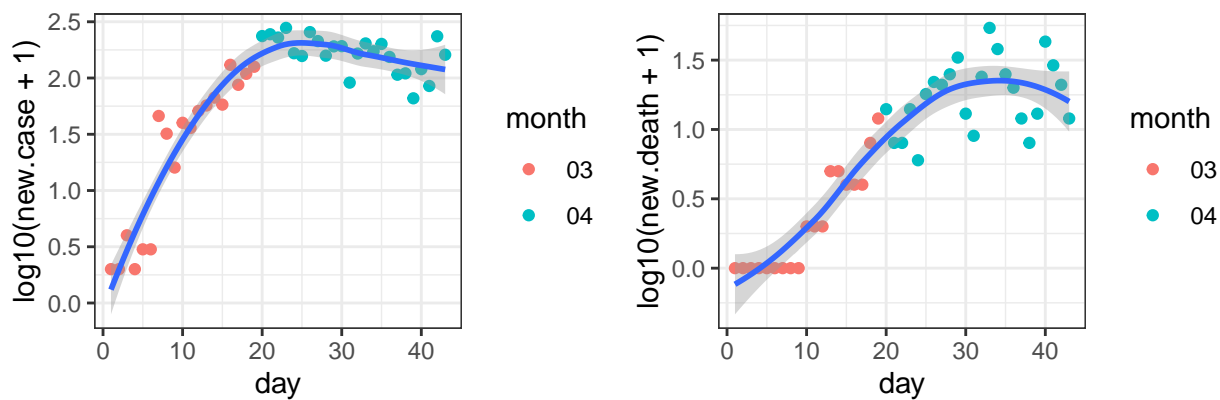
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-09

### Hartford\_Connecticut



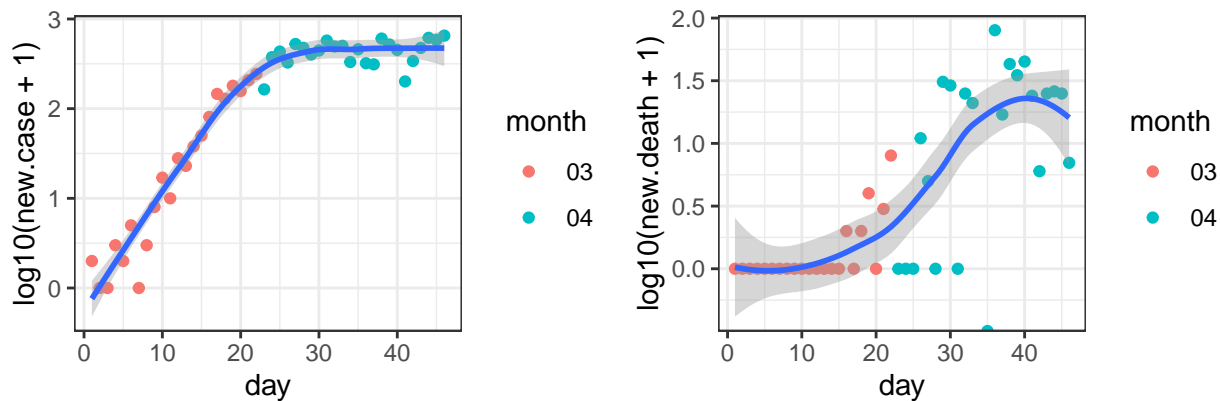
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-14

### Macomb\_Michigan



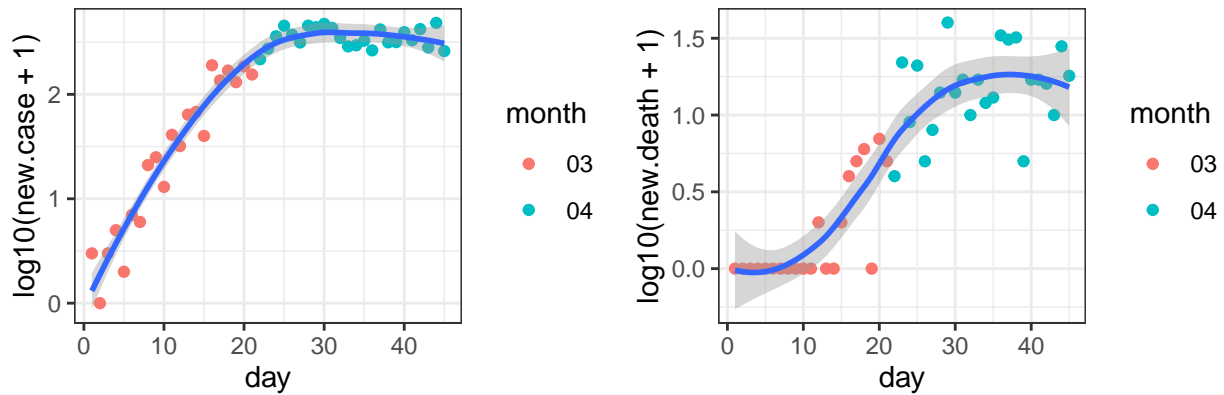
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-13

### Philadelphia\_Pennsylvania



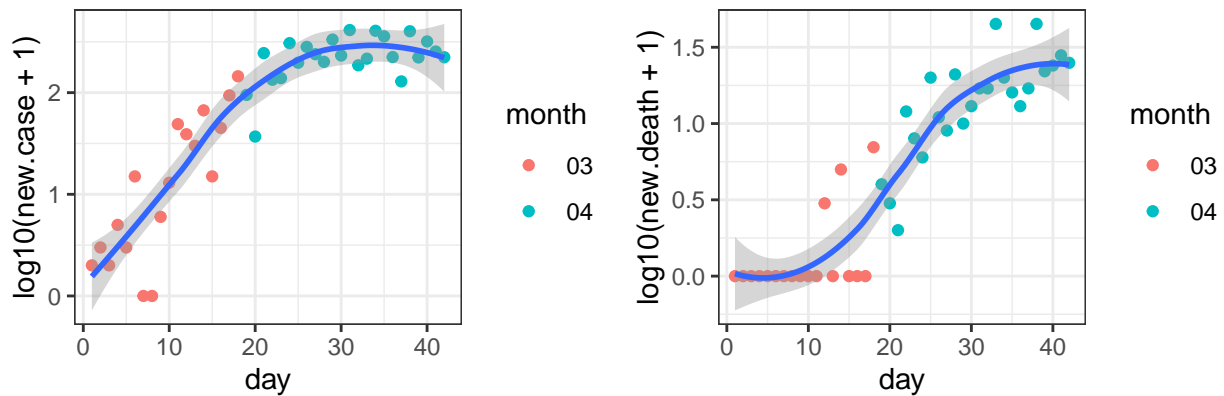
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-10

### Middlesex\_New Jersey



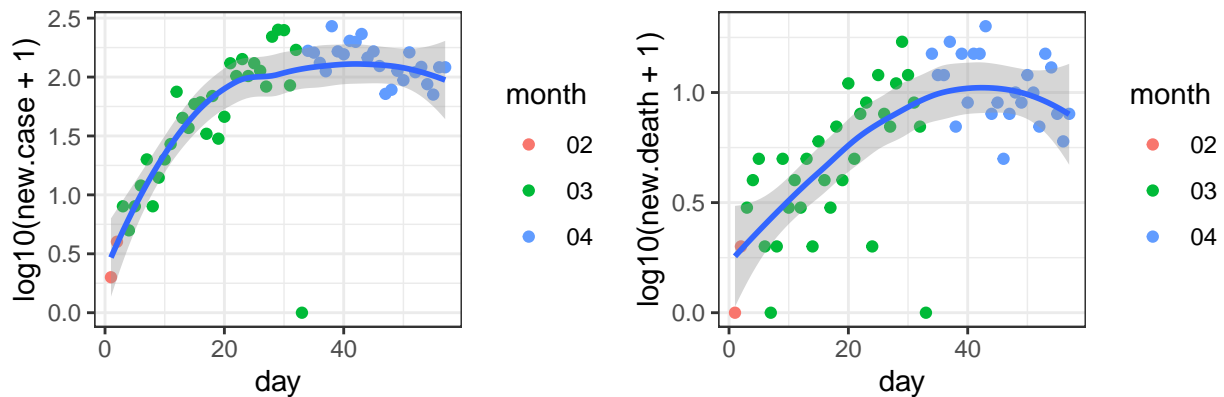
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-11

### New Haven\_Connecticut



data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-14

### King\_Washington



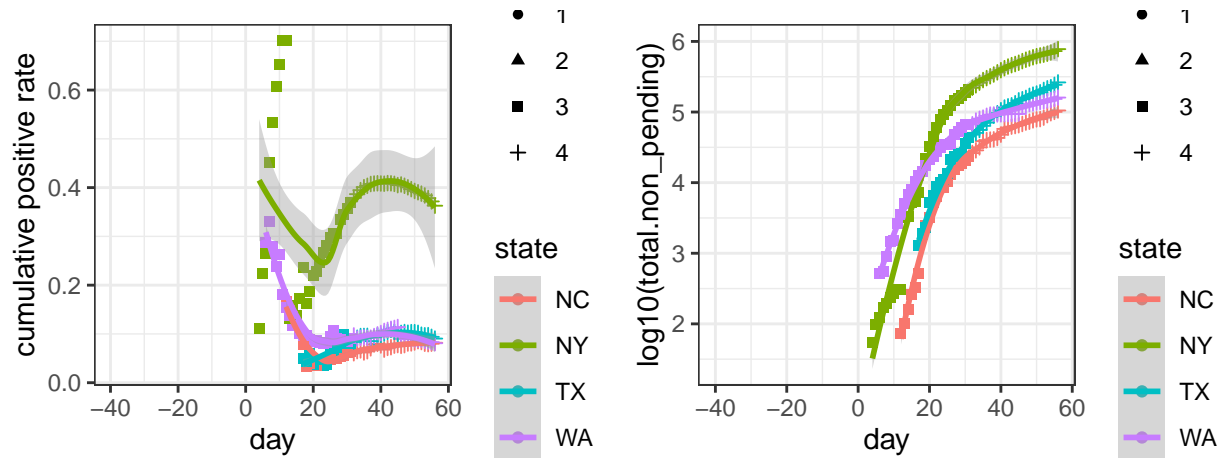
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 02-28

## COVID Tracking

The positive rates of testing can be an indicator on how much the COVID-19 has spread. However, they are more noisy data since the negative testing results are often not reported and the tests are almost surely taken on a non-representative random sample of the population. The COVID tracking project provides a grade per state: "If you are calculating positive rates, it should only be with states that have an A grade. And be

careful going back in time because almost all the states have changed their level of reporting at different times.” (<https://covidtracking.com/about-tracker/>). The data are also available for both counties and states, here I only look at state level data.

Since the daily positive rate can fluctuate a lot, here I only illustrate the cumulative positive rate across time, for four states with grade A data. Of course since this is an R markdown file, you can modify the source code and check for other states.



[github.com/COVID19Tracking/](https://github.com/COVID19Tracking/), cumulative positive rate on 0425: 0.08(WA) 0.09(TX) 0.36(NY) 0.08(NC)

## Session information

```
sessionInfo()
```

```
## R version 3.6.2 (2019-12-12)
## Platform: x86_64-apple-darwin15.6.0 (64-bit)
## Running under: macOS Catalina 10.15.4
##
## Matrix products: default
## BLAS: /Library/Frameworks/R.framework/Versions/3.6/Resources/lib/libRblas.0.dylib
## LAPACK: /Library/Frameworks/R.framework/Versions/3.6/Resources/lib/libRlapack.dylib
##
## locale:
## [1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
##
## attached base packages:
## [1] stats      graphics  grDevices  utils      datasets  methods   base
##
## other attached packages:
## [1] httr_1.4.1      ggpubr_0.2.5  magrittr_1.5  ggplot2_3.2.1
##
## loaded via a namespace (and not attached):
## [1] Rcpp_1.0.3      pillar_1.4.3    compiler_3.6.2  tools_3.6.2
## [5] digest_0.6.23   evaluate_0.14    lifecycle_0.1.0  tibble_2.1.3
## [9] gtable_0.3.0    pkgconfig_2.0.3  rlang_0.4.4      yaml_2.2.1
## [13] xfun_0.12       gridExtra_2.3    withr_2.1.2      dplyr_0.8.4
## [17] stringr_1.4.0   knitr_1.28       grid_3.6.2       tidyselect_1.0.0
## [21] cowplot_1.0.0   glue_1.3.1       R6_2.4.1         rmarkdown_2.1
## [25] purrr_0.3.3     farver_2.0.3     scales_1.1.0     htmltools_0.4.0
```

```
## [29] assertthat_0.2.1 colorspace_1.4-1 ggsignif_0.6.0 labeling_0.3
## [33] stringi_1.4.5 lazyeval_0.2.2 munsell_0.5.0 crayon_1.3.4
```