

# Exploration of COVID-19 tracking data from multiple resources

Wei Sun

2020-04-14

## 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>

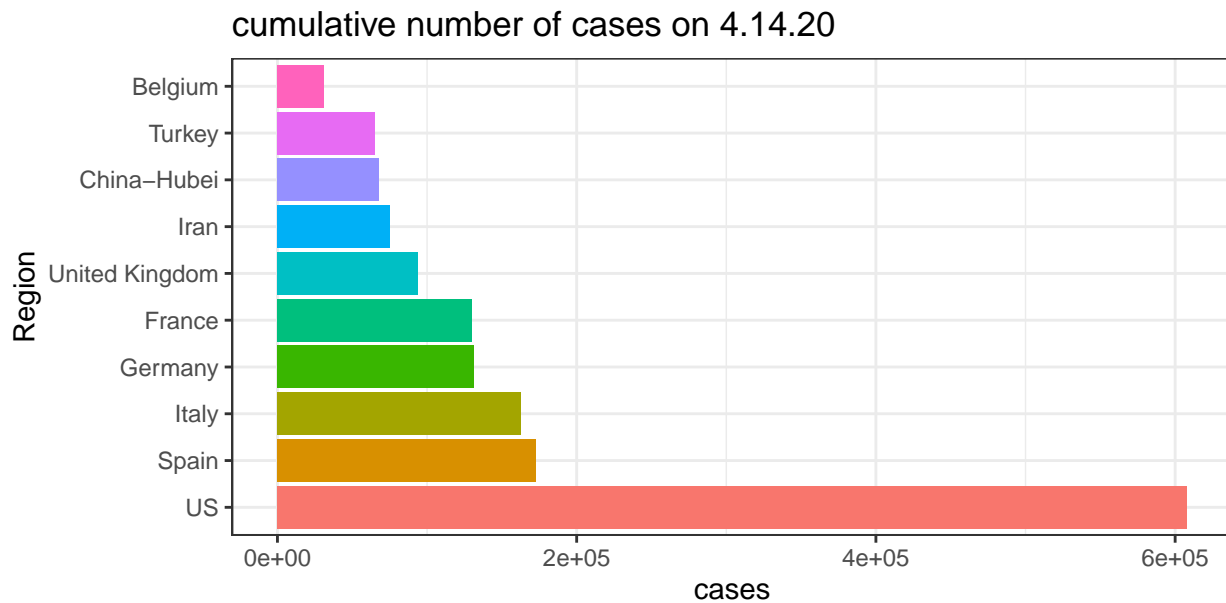
## JHU

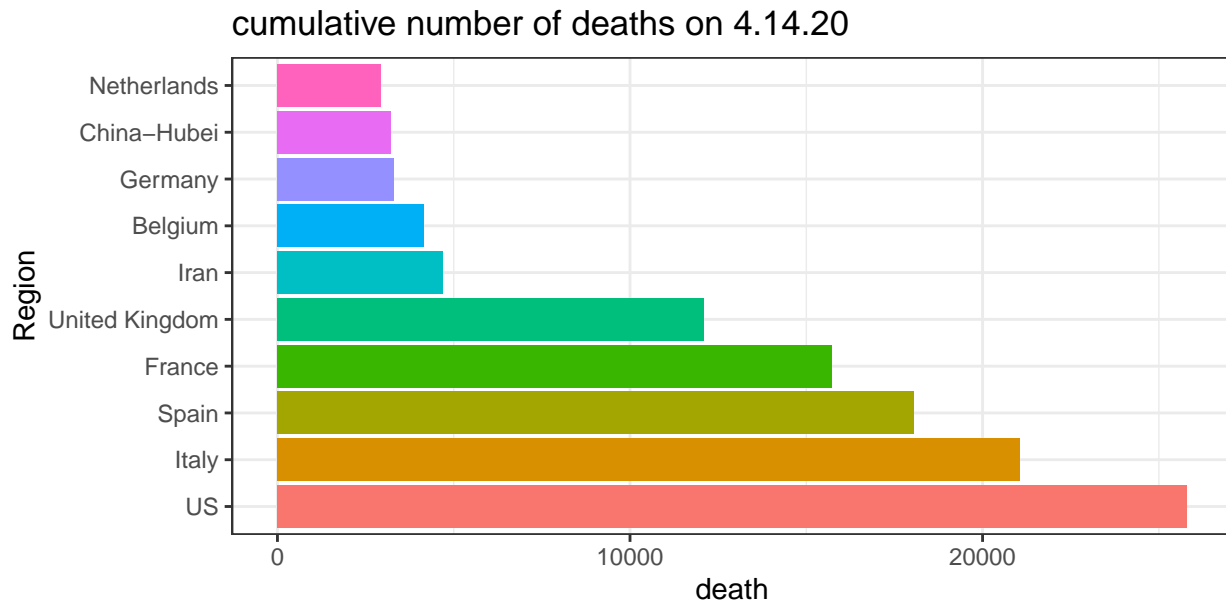
Assume you have cloned the JHU Github repository on your local machine at “../COVID-19”.

### 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.

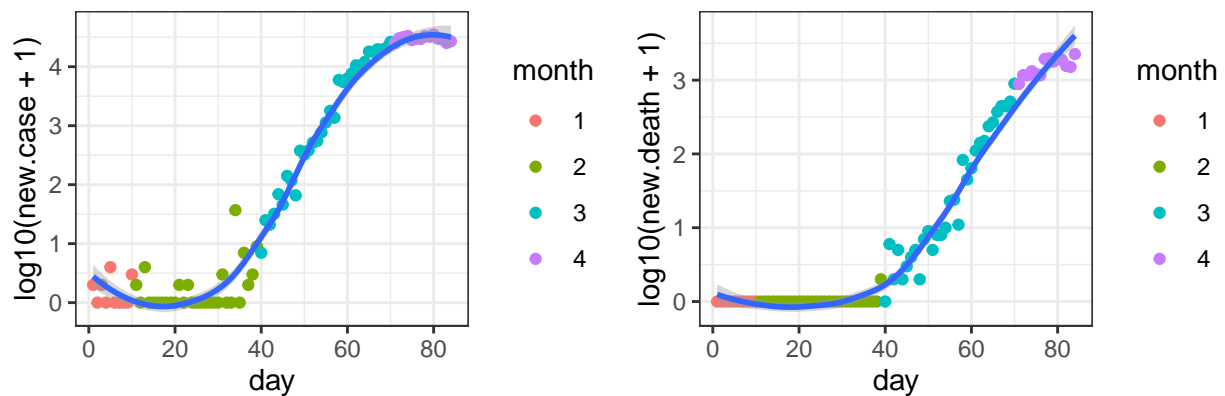
Here is the list of 10 records with the largest number of cases or deaths on the most recent date.





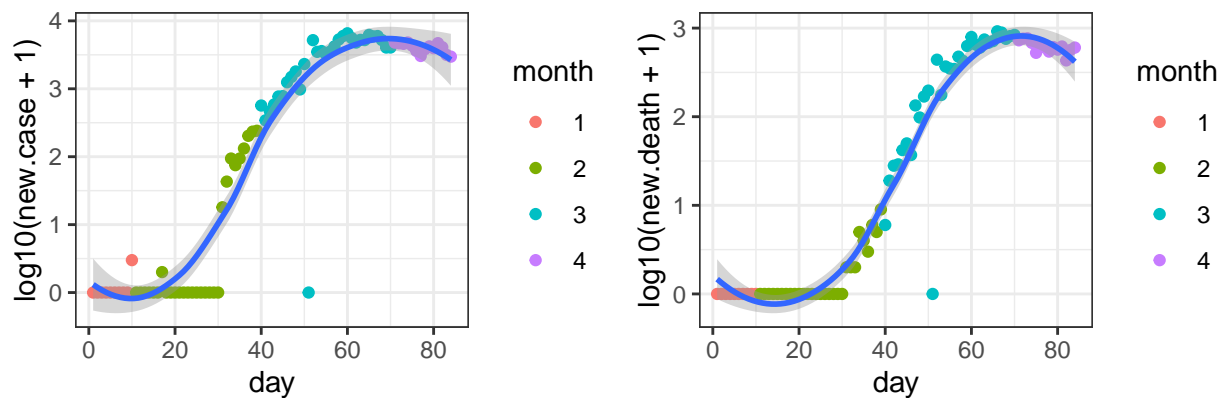
Next, I check for each country/region, what is the number of new cases/deaths? This data is important to understand what is the trend under different situations, e.g., population density, social distance policies etc. Here I checked the top 10 countries/regions with the highest number of deaths.

### US



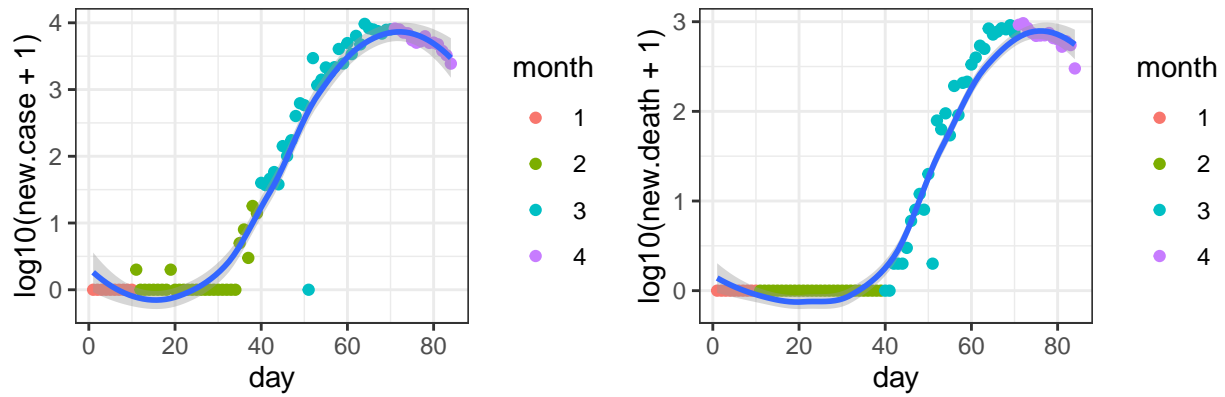
data source: <https://github.com/CSSEGISandData/COVID-19>, day 1 is 1/22/2020

### Italy



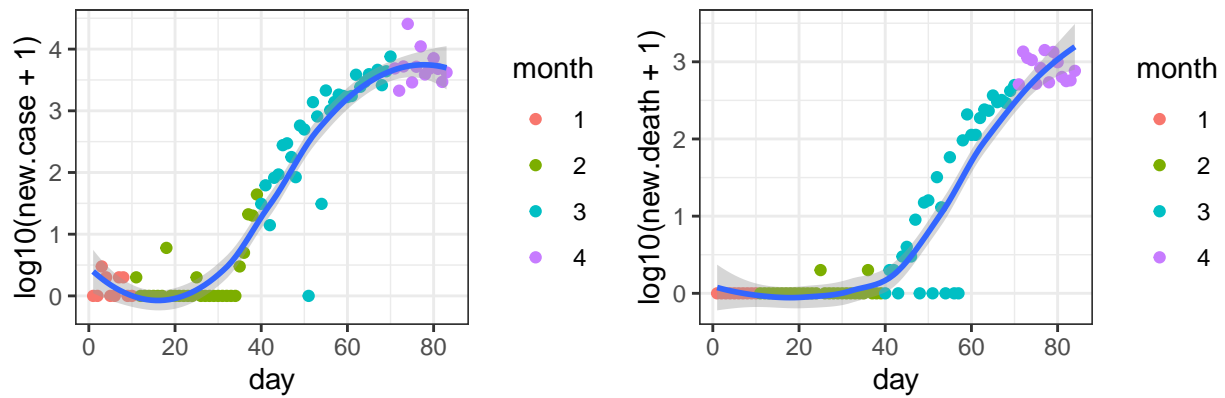
data source: <https://github.com/CSSEGISandData/COVID-19>, day 1 is 1/22/2020

### Spain



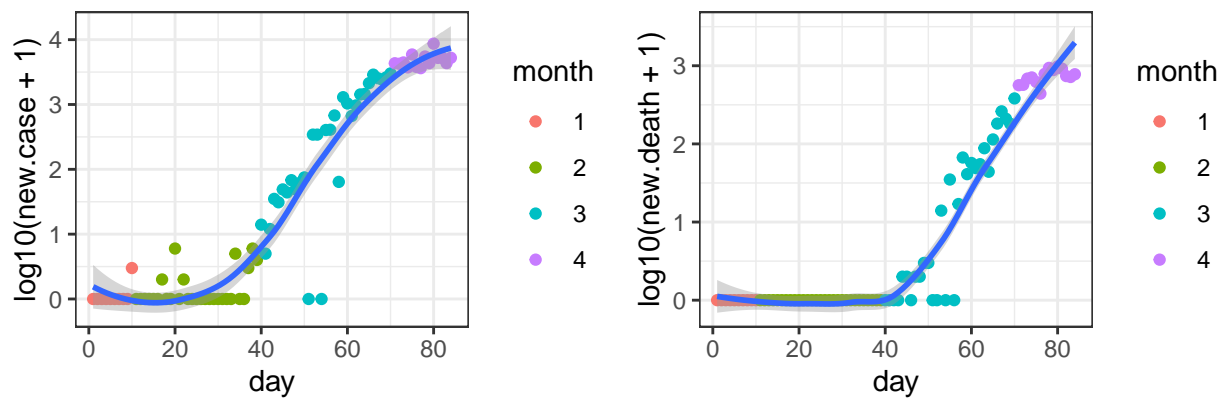
data source: <https://github.com/CSSEGISandData/COVID-19>, day 1 is 1/22/2020

### France



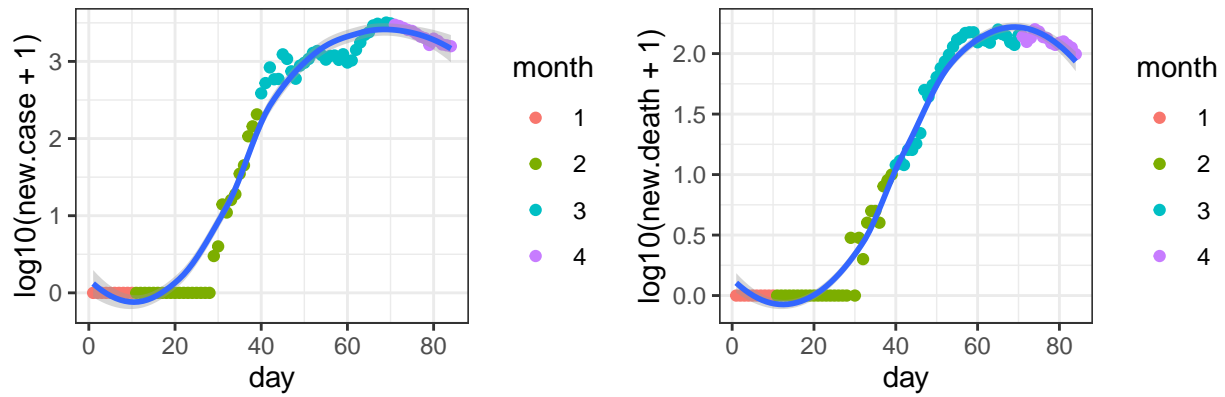
data source: <https://github.com/CSSEGISandData/COVID-19>, day 1 is 1/22/2020

### United Kingdom



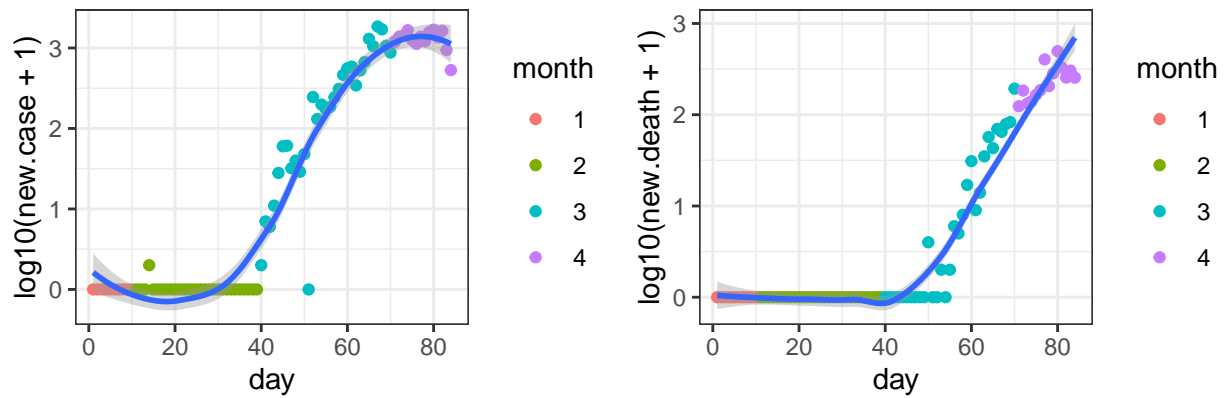
data source: <https://github.com/CSSEGISandData/COVID-19>, day 1 is 1/22/2020

### Iran



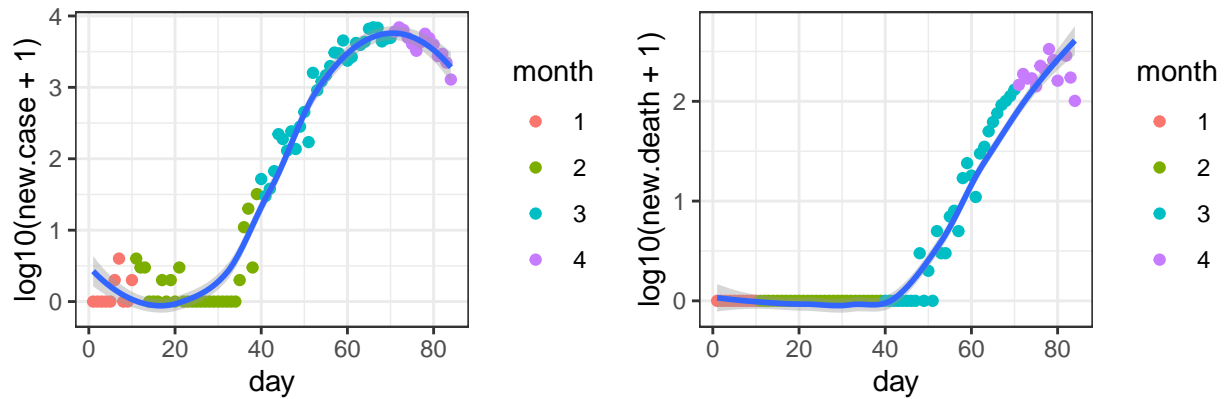
data source: <https://github.com/CSSEGISandData/COVID-19>, day 1 is 1/22/2020

### Belgium

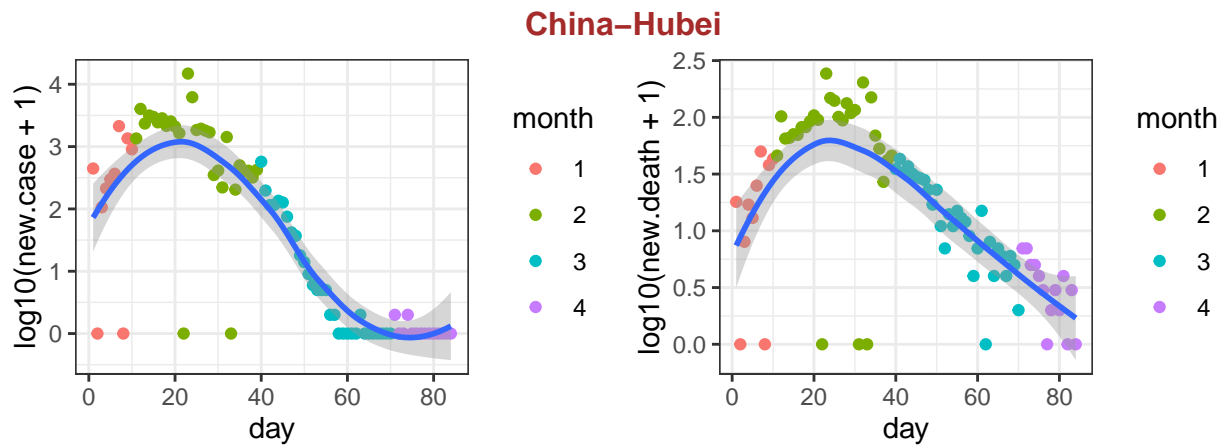


data source: <https://github.com/CSSEGISandData/COVID-19>, day 1 is 1/22/2020

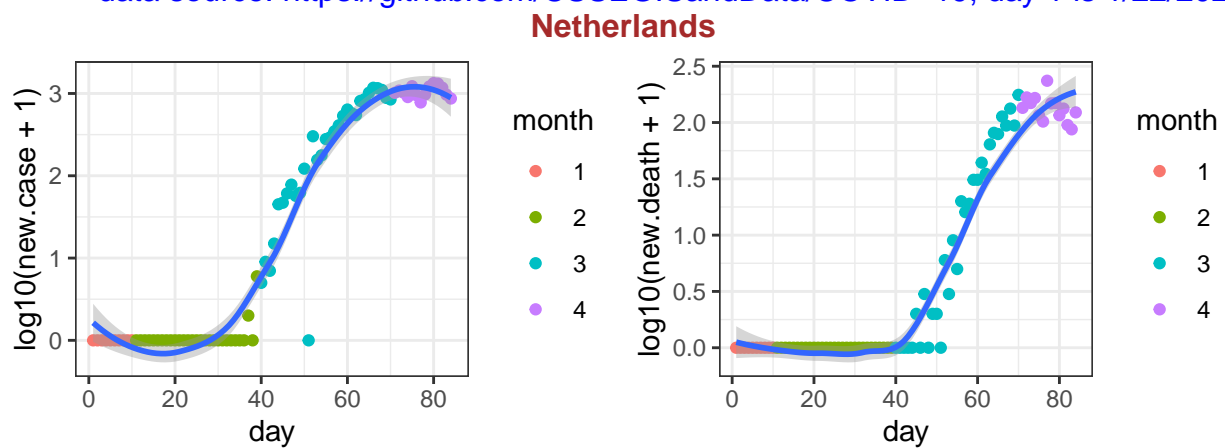
### Germany



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



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

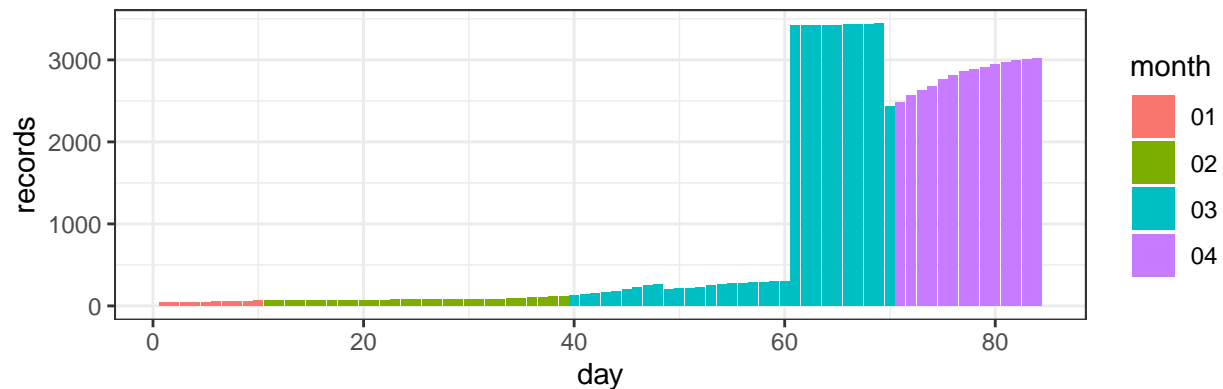


data source: <https://github.com/CSSEGISandData/COVID-19>, day 1 is 1/22/2020

## 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

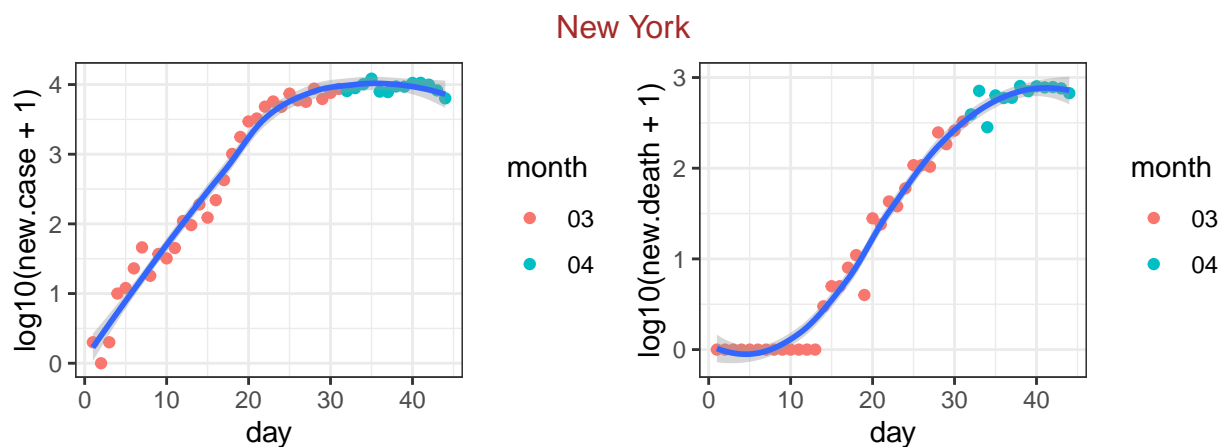
```
## [1] "2020-04-13"
```

### state level data

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

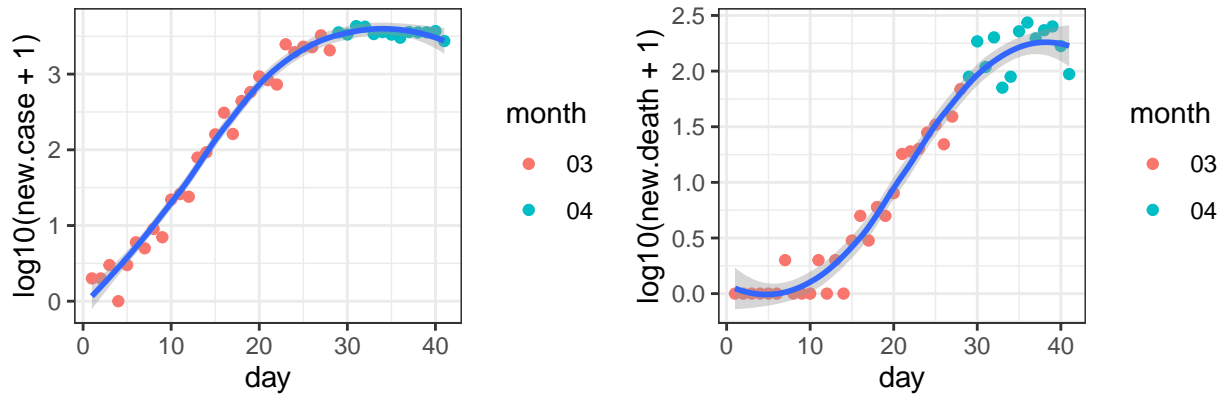
##	date	state	fips	cases	deaths
## 2308	2020-04-13	New York	36	195031	10056
## 2306	2020-04-13	New Jersey	34	64584	2443
## 2298	2020-04-13	Michigan	26	25487	1601
## 2294	2020-04-13	Louisiana	22	21016	884
## 2297	2020-04-13	Massachusetts	25	26867	844
## 2289	2020-04-13	Illinois	17	22025	800
## 2279	2020-04-13	California	6	24334	725
## 2281	2020-04-13	Connecticut	9	13381	602
## 2315	2020-04-13	Pennsylvania	42	24295	563
## 2326	2020-04-13	Washington	53	10538	525
## 2284	2020-04-13	Florida	12	21011	498
## 2285	2020-04-13	Georgia	13	13125	479
## 2290	2020-04-13	Indiana	18	8236	350
## 2321	2020-04-13	Texas	48	14488	320
## 2280	2020-04-13	Colorado	8	7691	308
## 2312	2020-04-13	Ohio	39	6975	274
## 2296	2020-04-13	Maryland	24	8936	262
## 2328	2020-04-13	Wisconsin	55	3428	155
## 2325	2020-04-13	Virginia	51	5747	149
## 2301	2020-04-13	Missouri	29	4388	137

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.



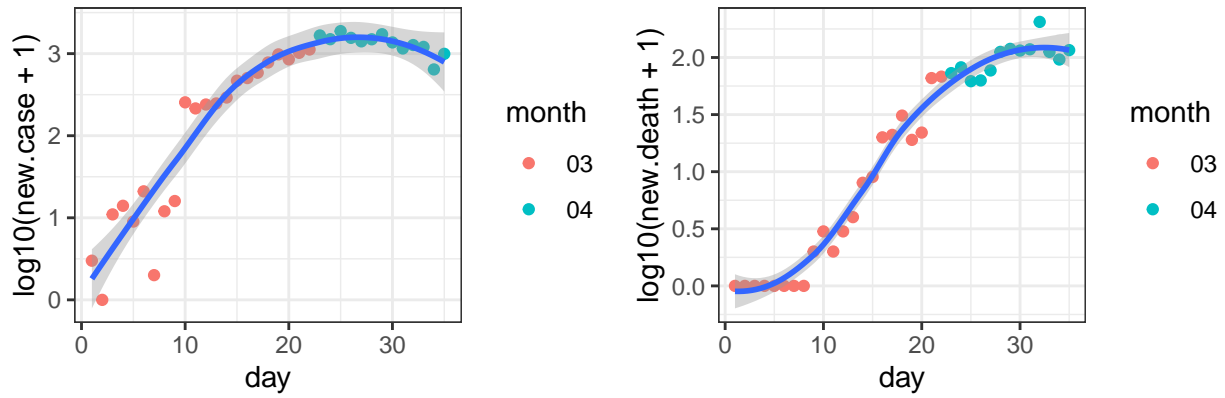
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-01

### New Jersey



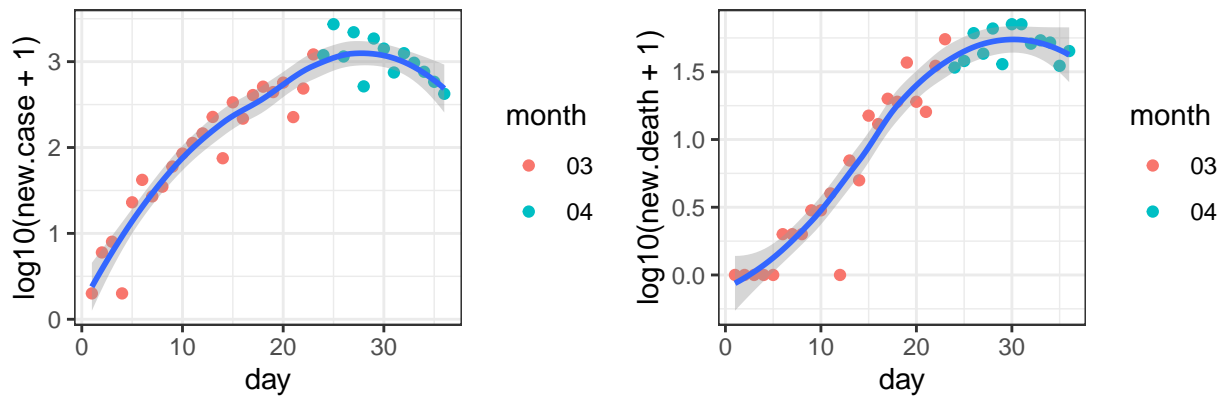
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-04

### Michigan



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

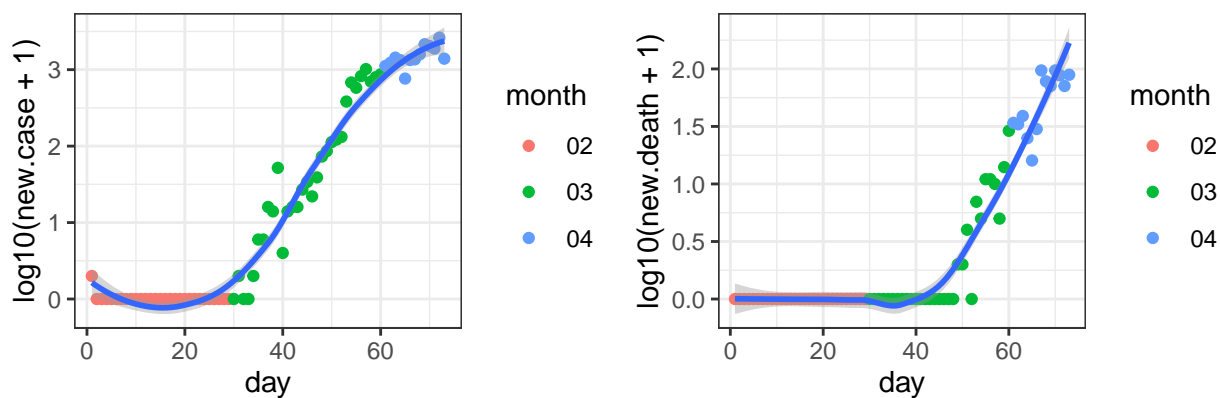
### Louisiana



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

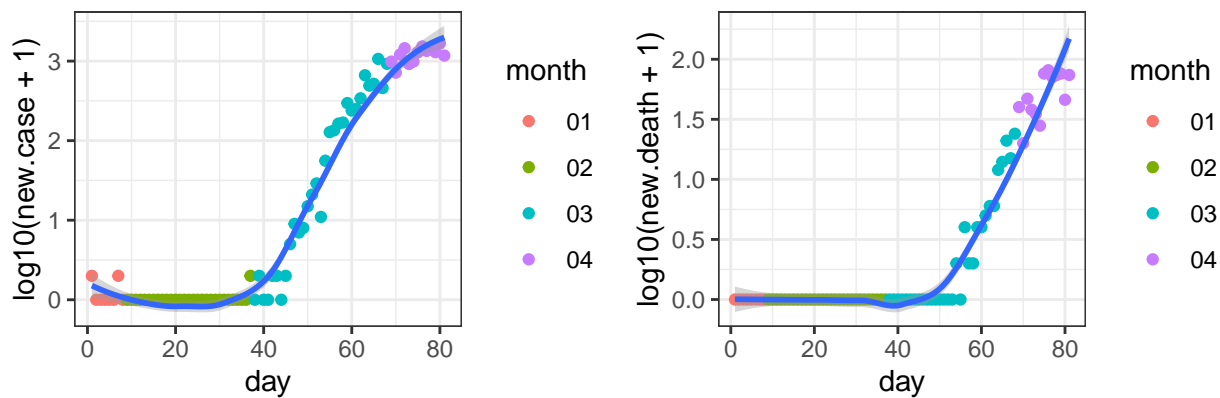


### Massachusetts



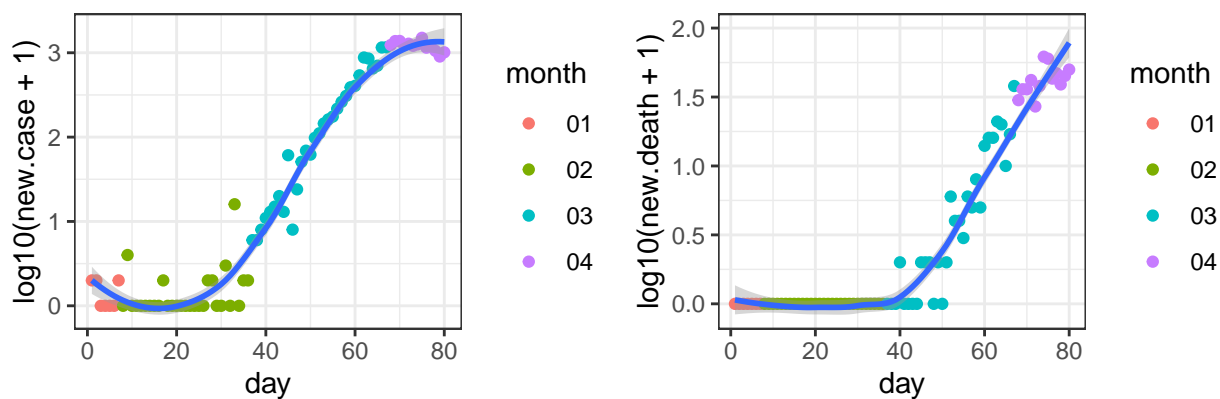
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 02-01

### Illinois



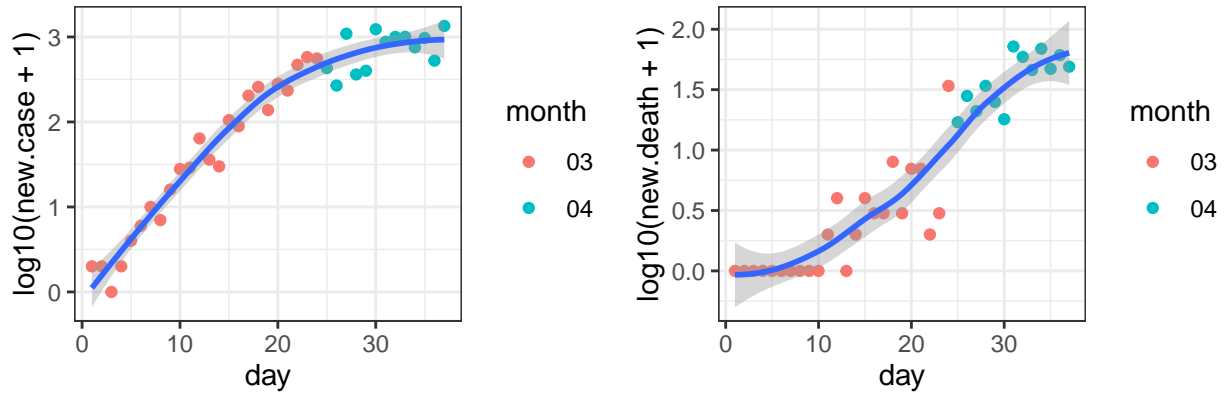
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 01-24

### California



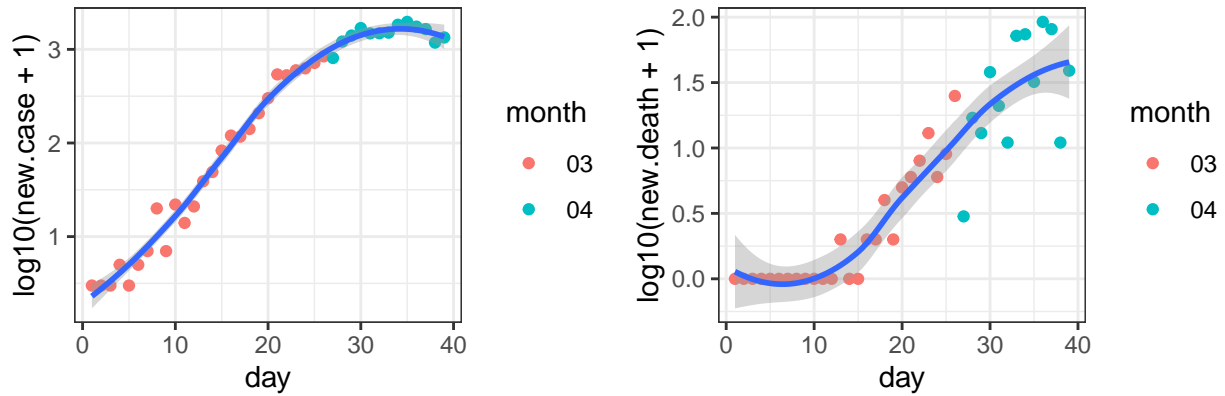
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 01-25

### Connecticut



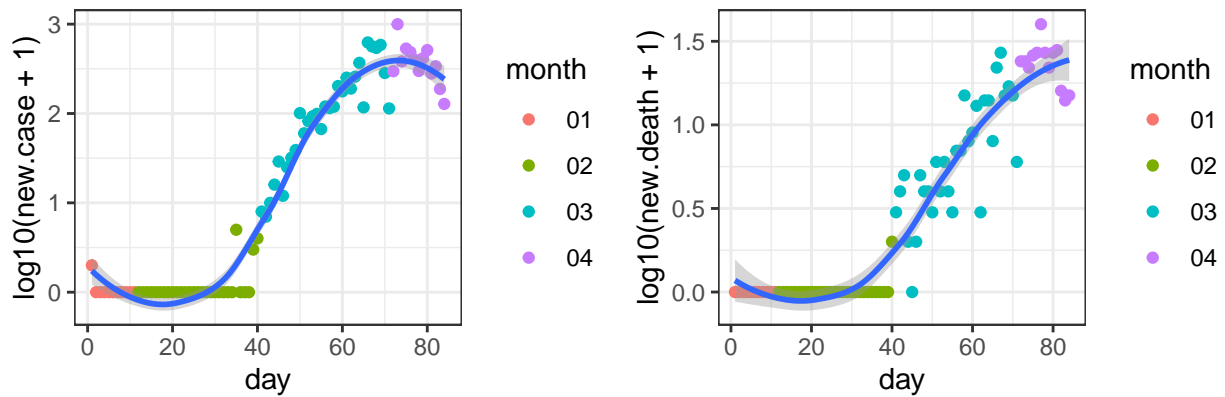
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-08

### Pennsylvania



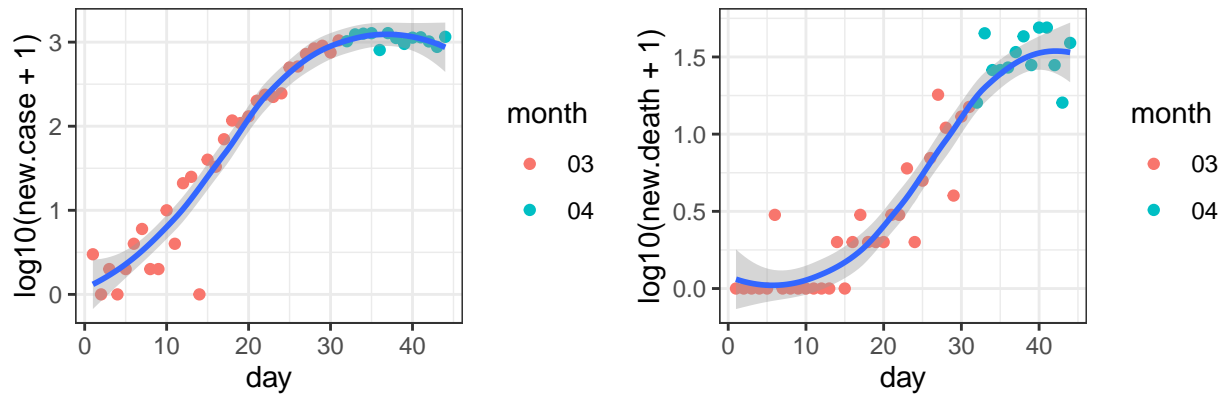
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-06

### Washington



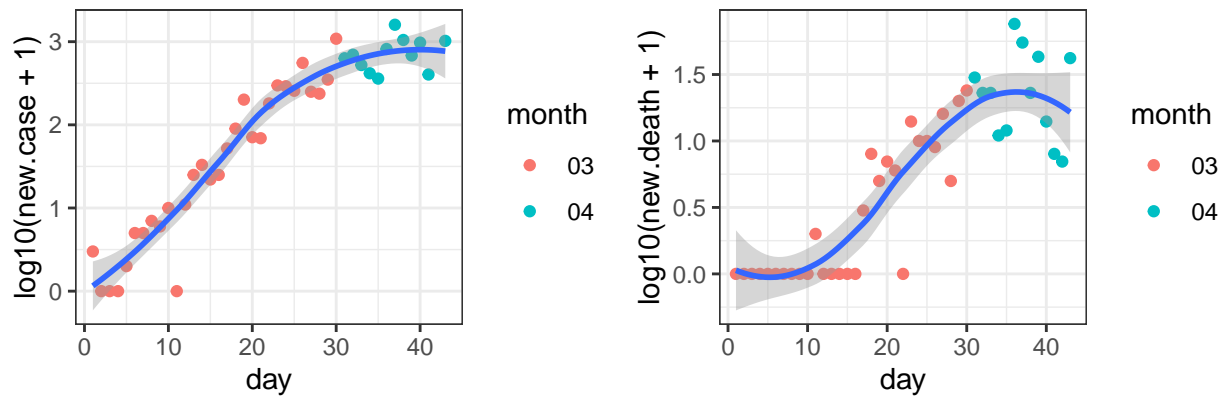
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 01-21

## Florida



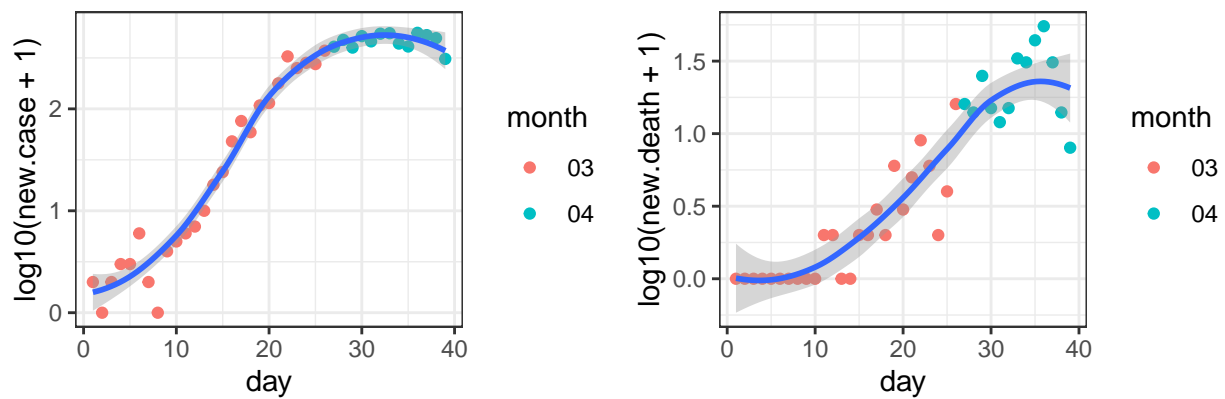
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-01

## Georgia



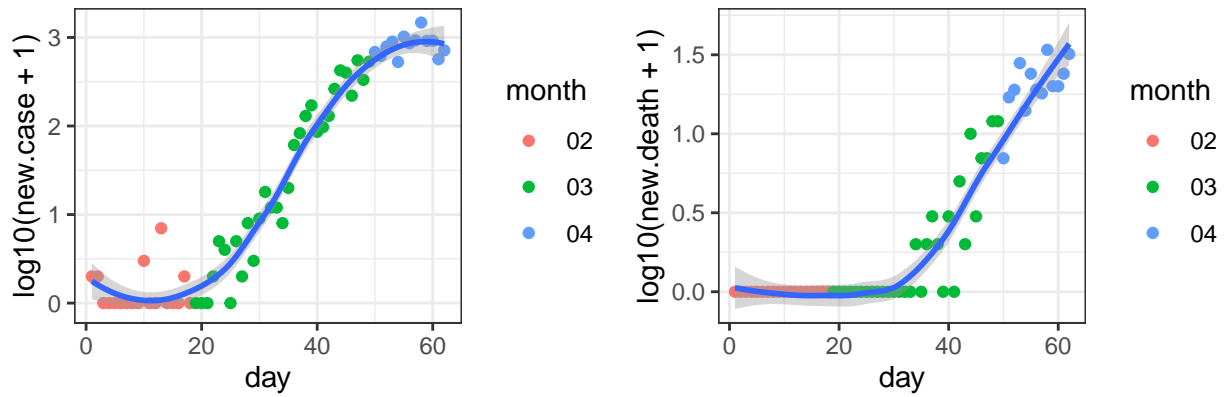
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-02

## Indiana



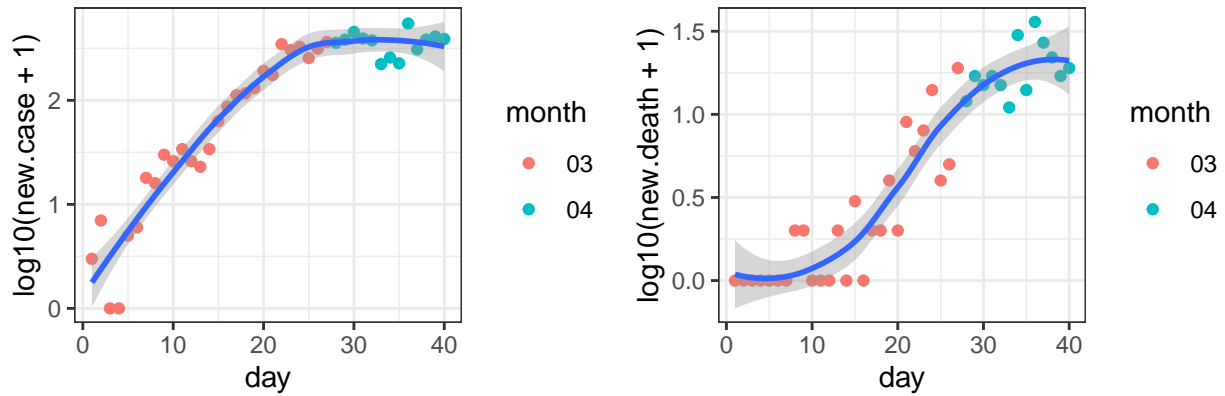
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-06

## Texas



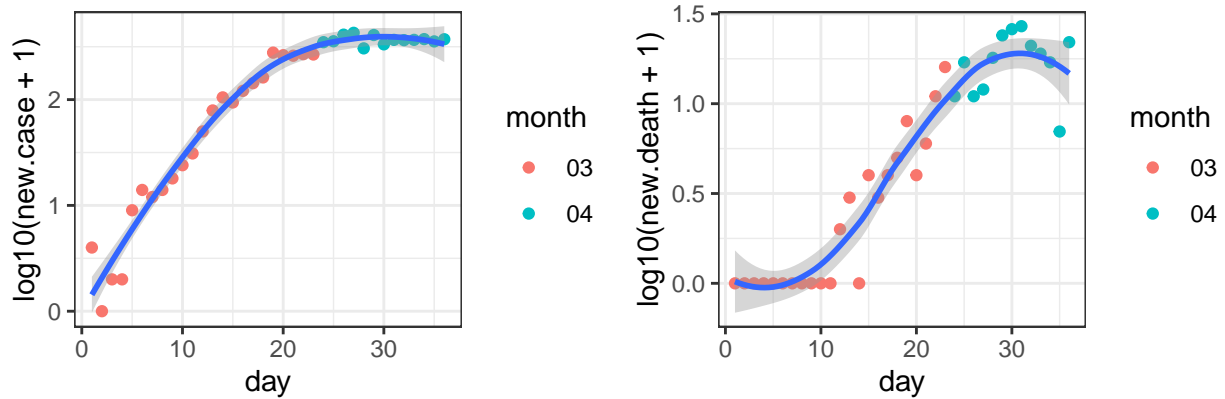
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 02-12

## Colorado



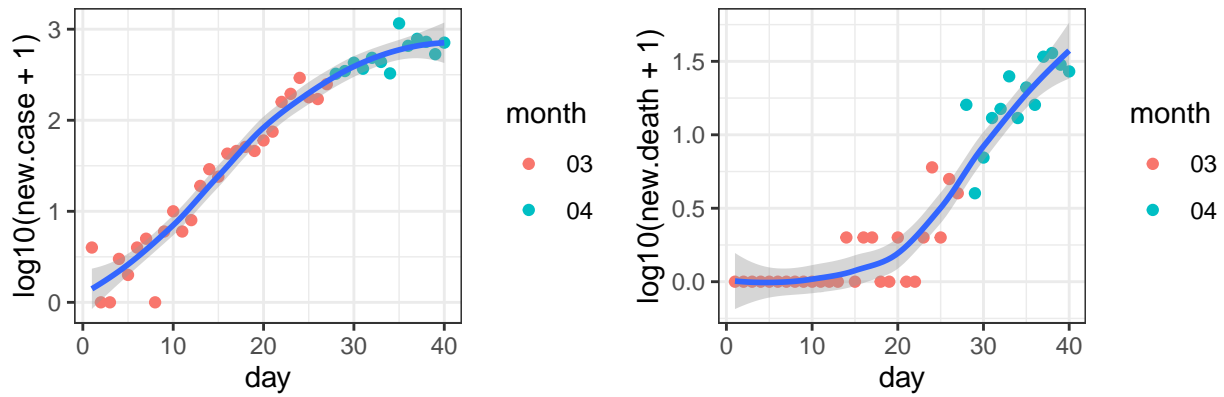
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-05

## Ohio



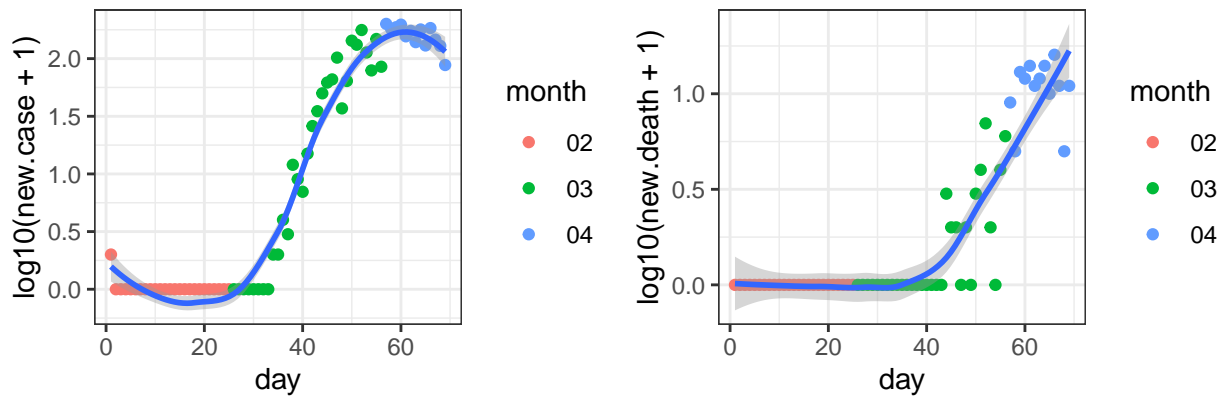
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-09

### Maryland



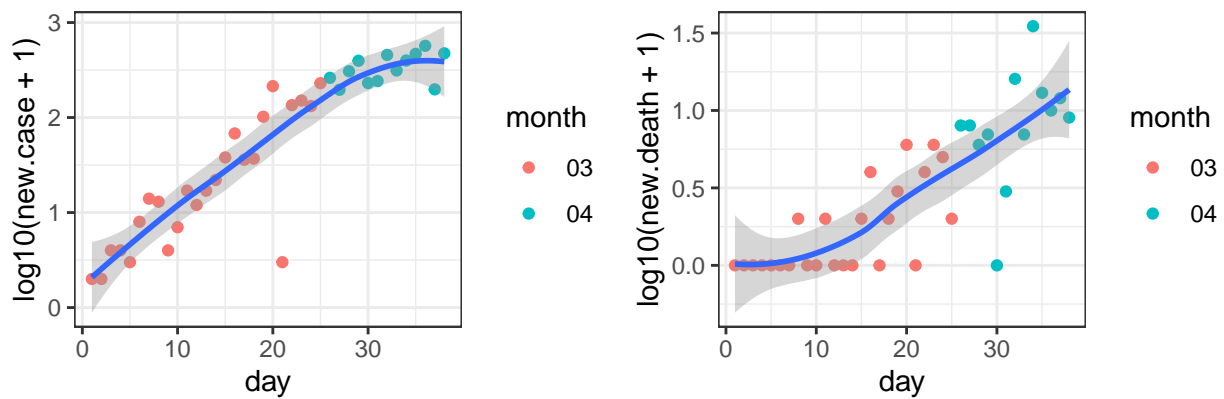
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-05

### Wisconsin

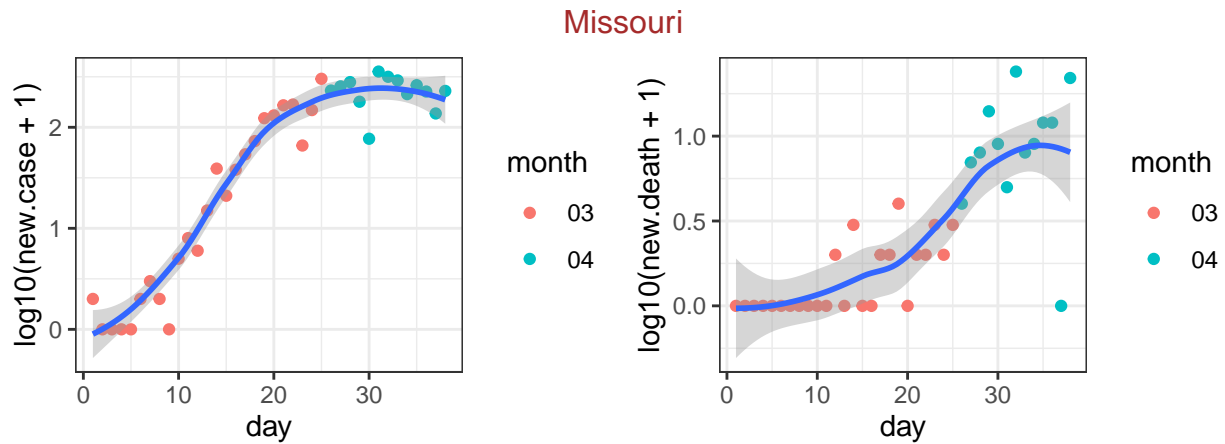


data source: <https://github.com/nytimes/covid-19-data>, day 1 is 02-05

### Virginia

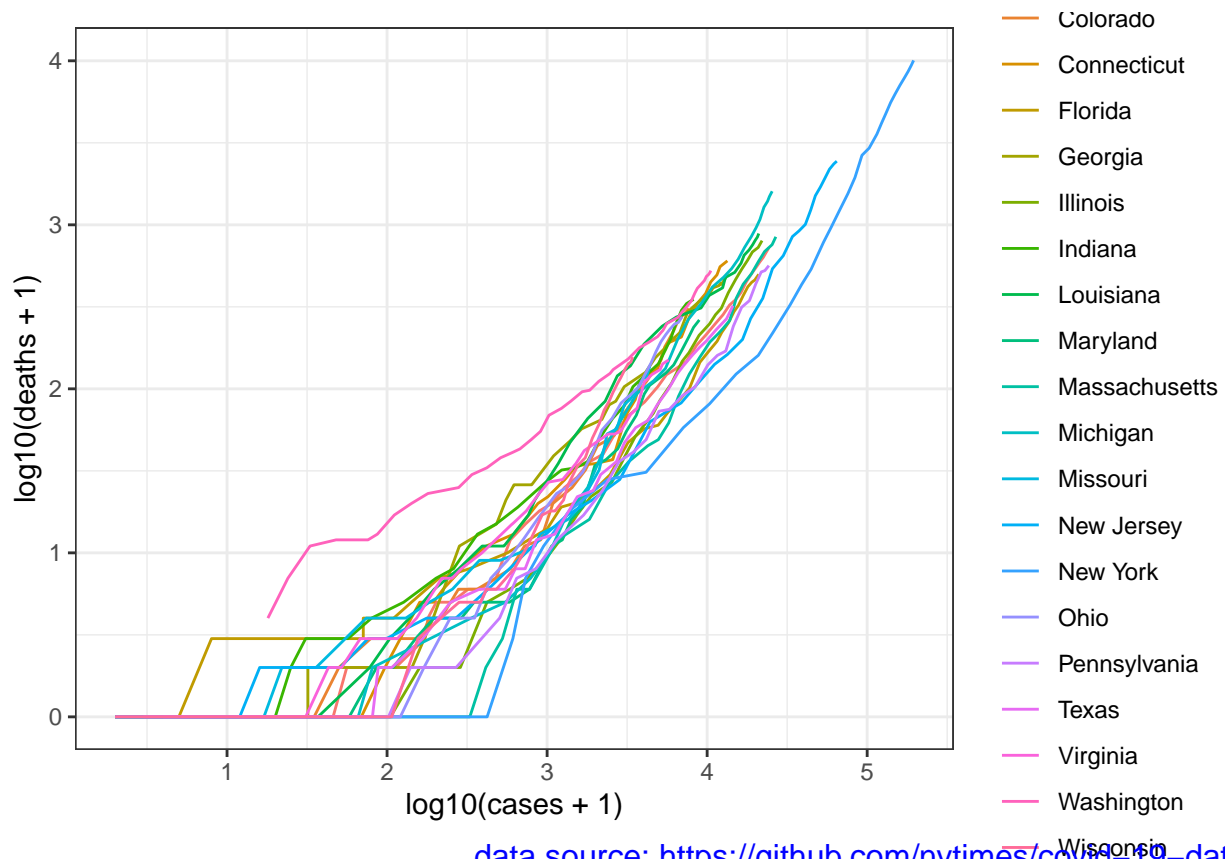


data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-07



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

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

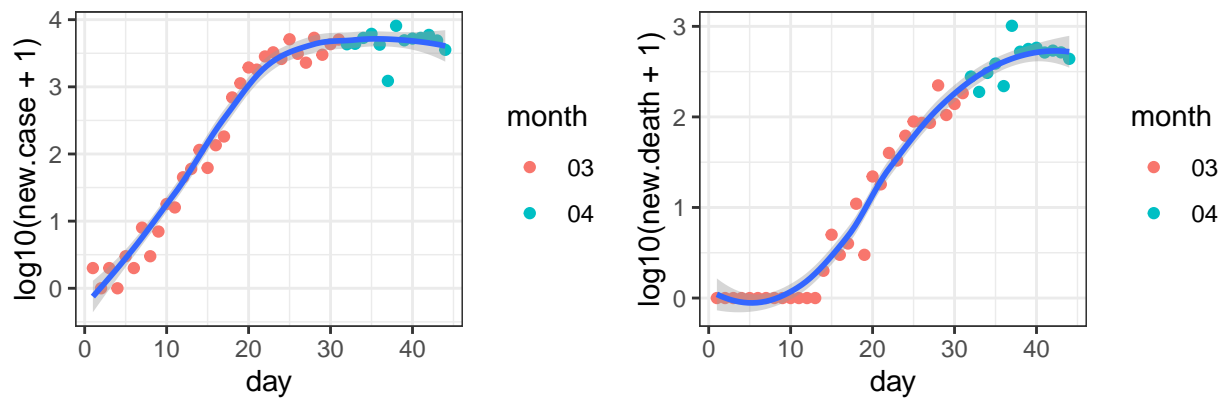
First check the 20 counties with the largest number of deaths.

##	date	county	state	fips	cases	deaths
## 55443	2020-04-13	New York City	New York	NA	106764	7154
## 55442	2020-04-13	Nassau	New York	36059	24358	1109
## 55027	2020-04-13	Wayne	Michigan	26163	11648	760

##	55470	2020-04-13	Westchester	New York	36119	19785	610
##	55462	2020-04-13	Suffolk	New York	36103	21643	580
##	54414	2020-04-13	Cook	Illinois	17031	15474	543
##	55368	2020-04-13	Bergen	New Jersey	34003	10092	482
##	55373	2020-04-13	Essex	New Jersey	34013	7634	433
##	55008	2020-04-13	Oakland	Michigan	26125	5073	347
##	54032	2020-04-13	Los Angeles	California	6037	9420	320
##	56386	2020-04-13	King	Washington	53033	4551	298
##	54125	2020-04-13	Fairfield	Connecticut	9001	6004	262
##	54867	2020-04-13	Orleans	Louisiana	22071	5651	244
##	54995	2020-04-13	Macomb	Michigan	26099	3418	240
##	55375	2020-04-13	Hudson	New Jersey	34017	7879	236
##	55386	2020-04-13	Union	New Jersey	34039	6636	217
##	55378	2020-04-13	Middlesex	New Jersey	34023	5987	204
##	54857	2020-04-13	Jefferson	Louisiana	22051	5088	186
##	55454	2020-04-13	Rockland	New York	36087	7965	182
##	54945	2020-04-13	Middlesex	Massachusetts	25017	5983	163

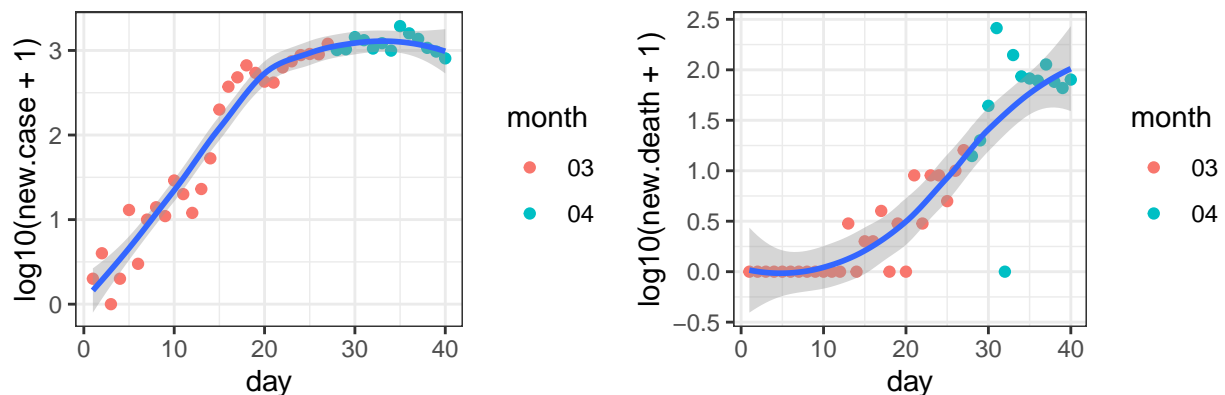
For these 20 counties, I check the number of new cases and the number of new deaths.

### New York City\_New York



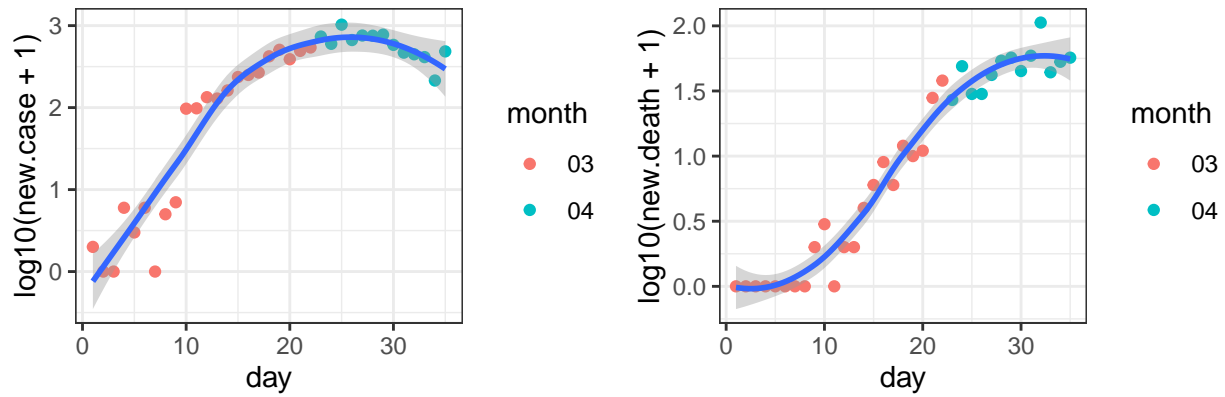
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-01

### Nassau\_New York



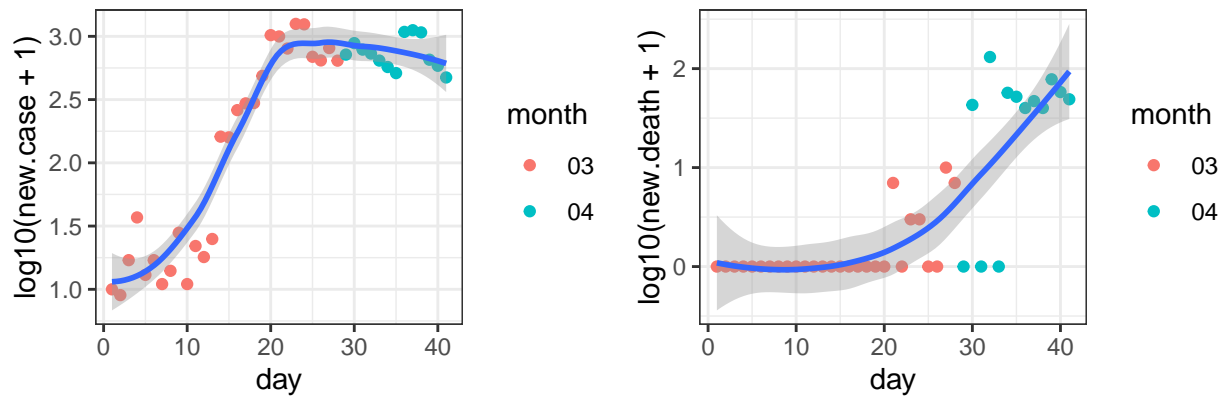
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-05

### Wayne\_Michigan



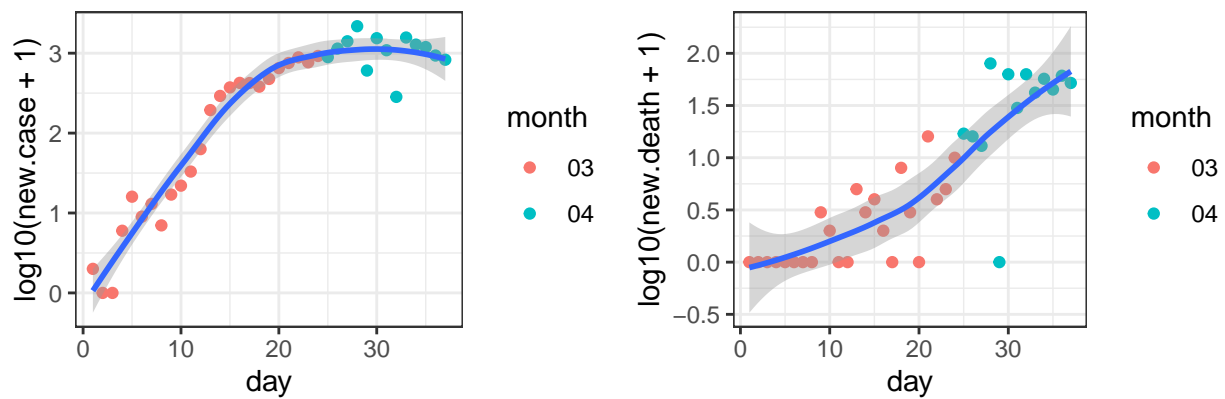
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-10

### Westchester\_New York



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

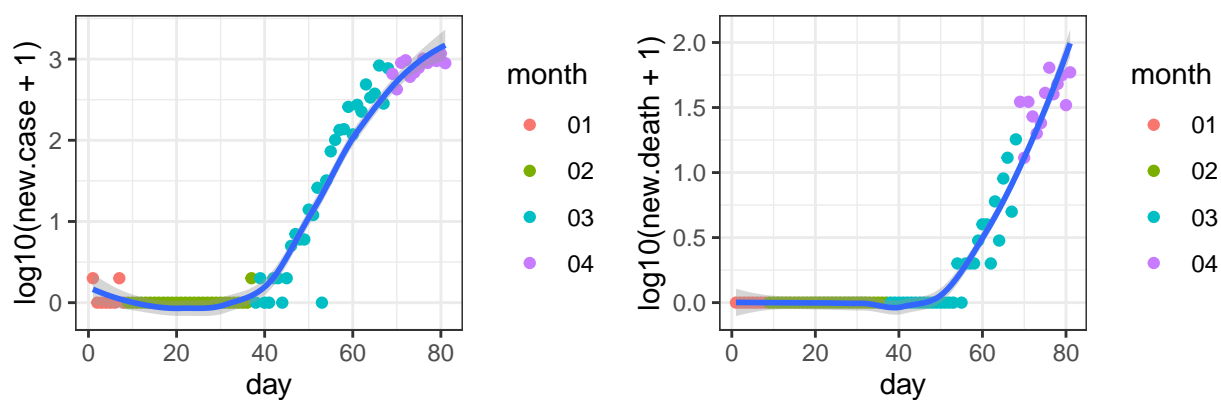
### Suffolk\_New York



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

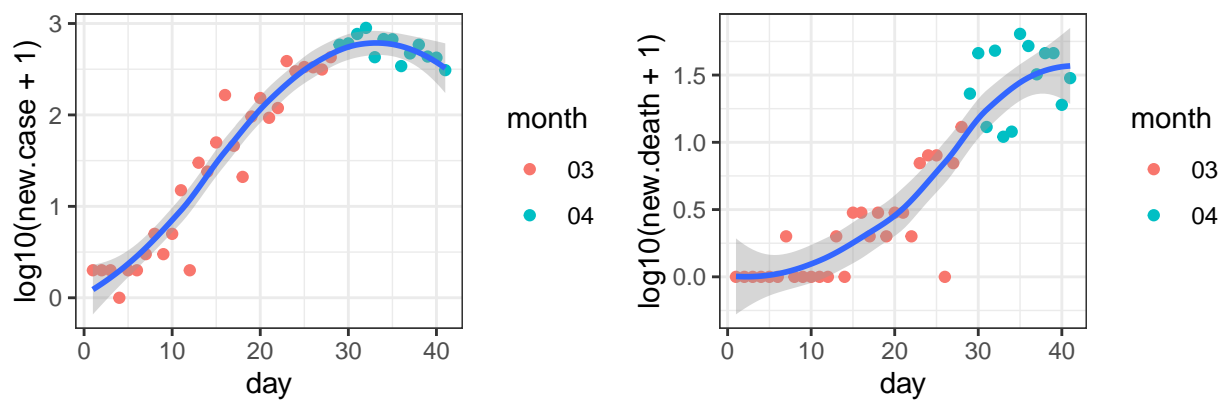


### Cook\_Illinois



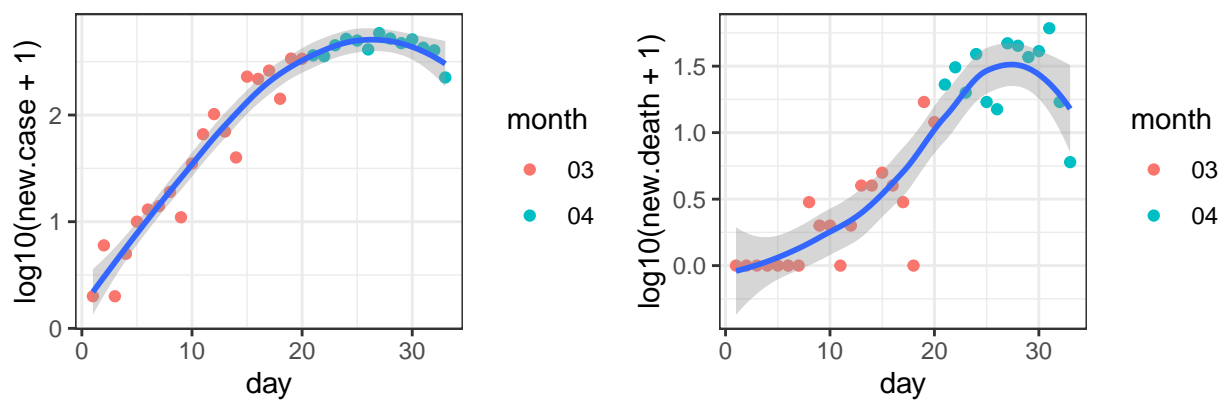
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 01-24

### Bergen\_New Jersey



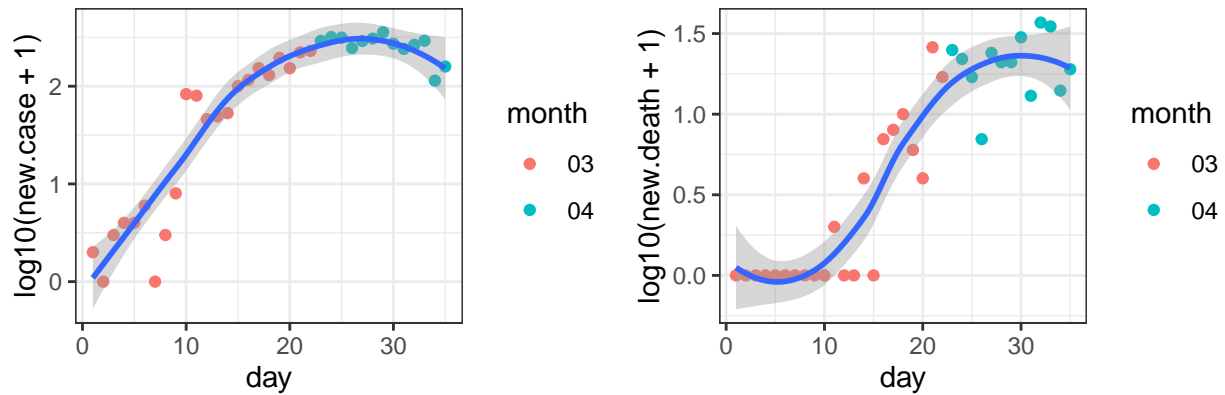
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-04

### Essex\_New Jersey



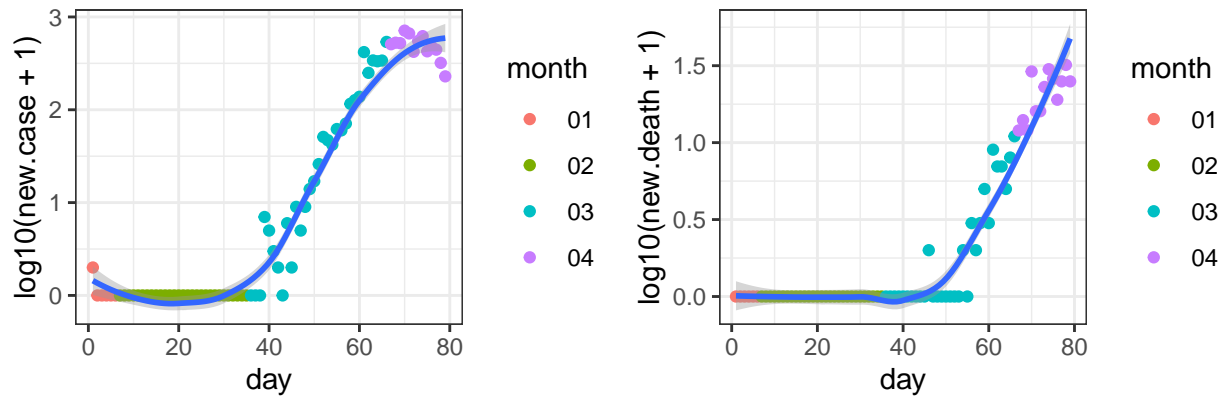
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-12

### Oakland\_Michigan



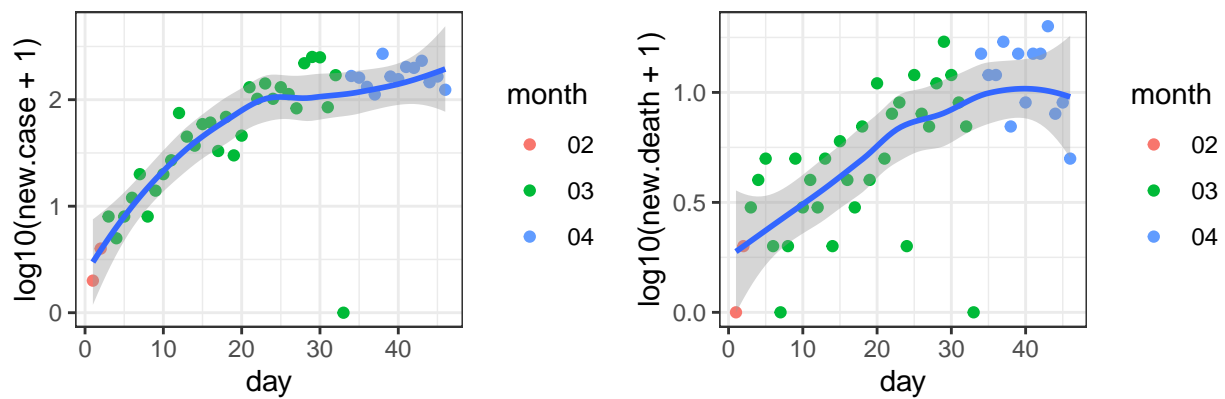
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-10

### Los Angeles\_California



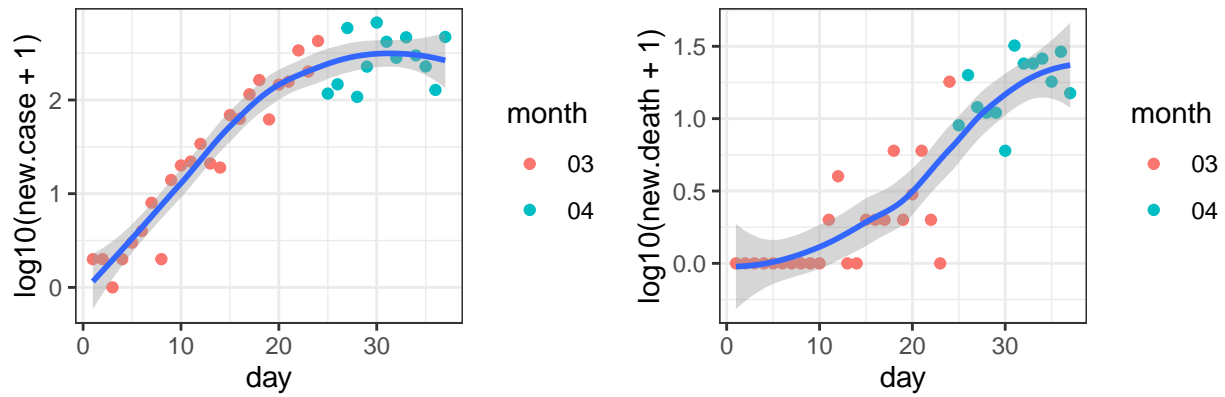
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 01-26

### King\_Washington



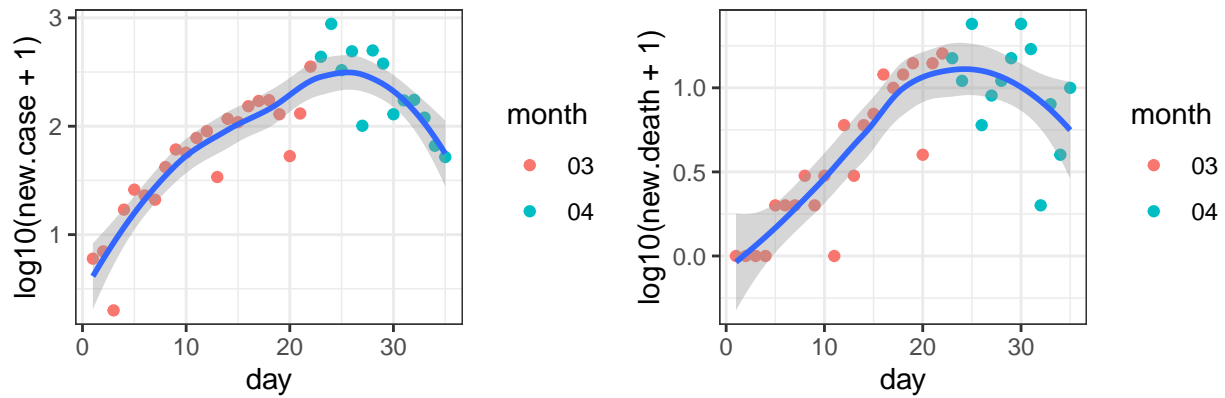
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 02-28

### Fairfield\_Connecticut



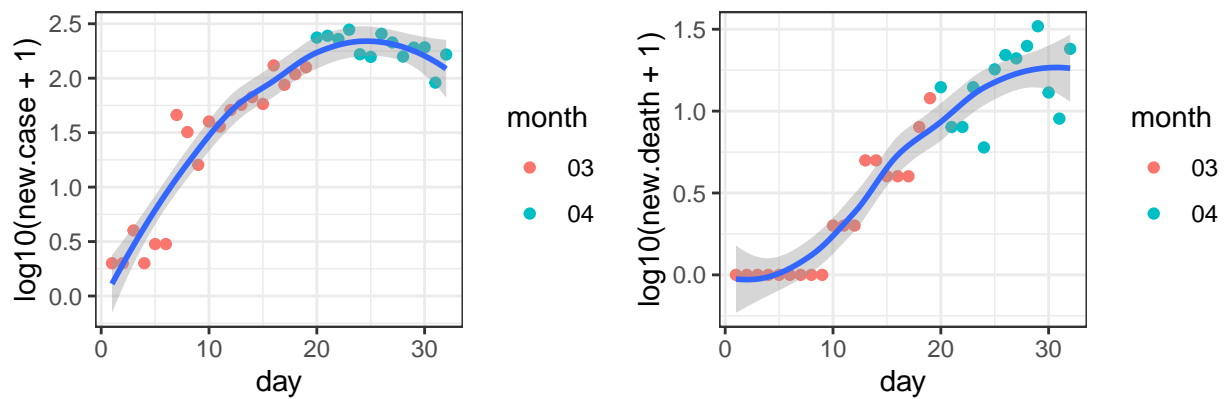
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-08

### Orleans\_Louisiana



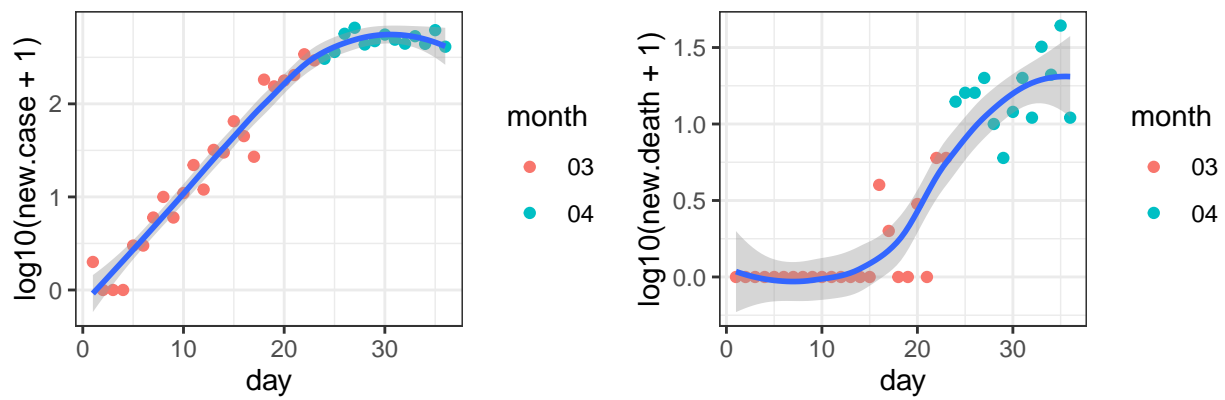
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-10

### Macomb\_Michigan



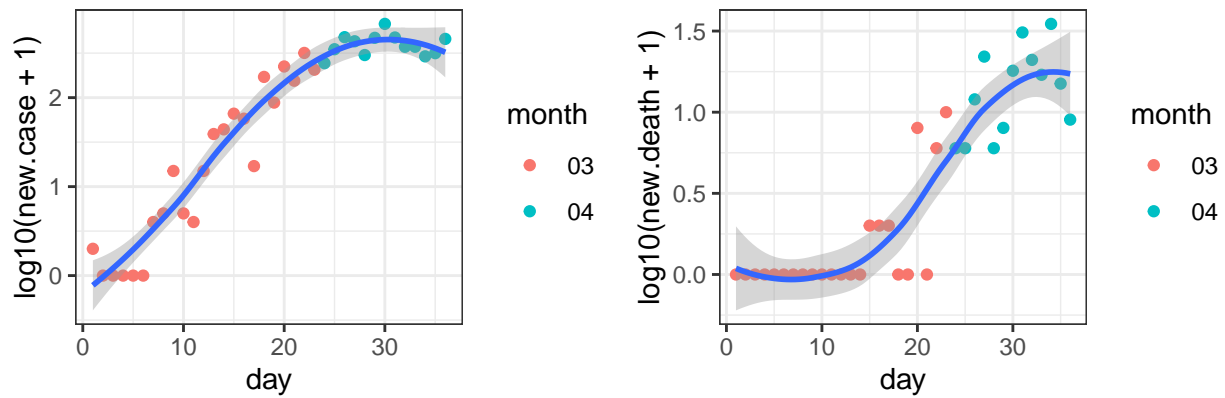
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-13

### Hudson\_New Jersey



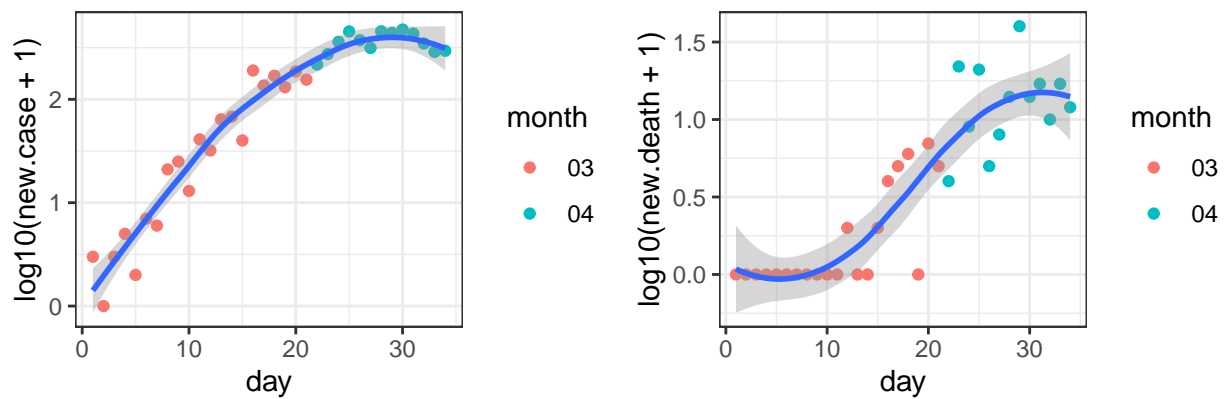
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-09

### Union\_New Jersey



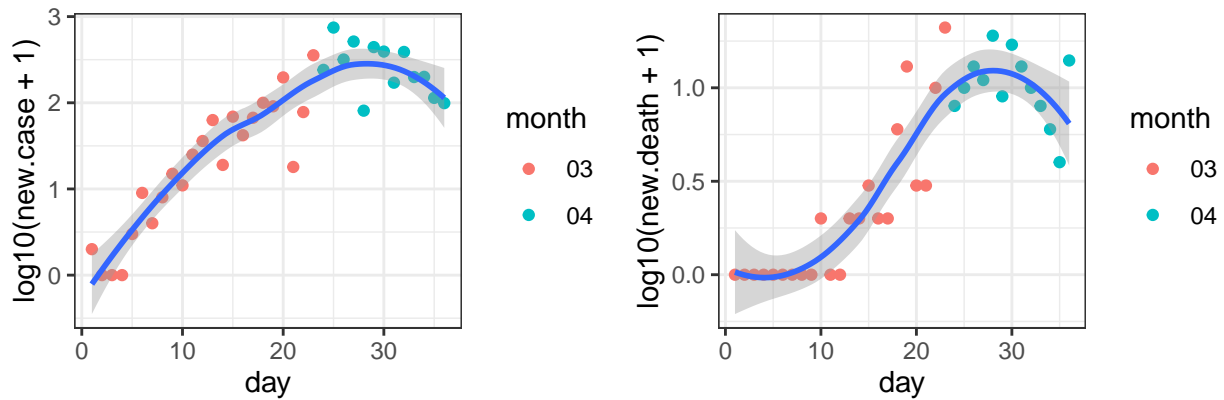
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-09

### Middlesex\_New Jersey



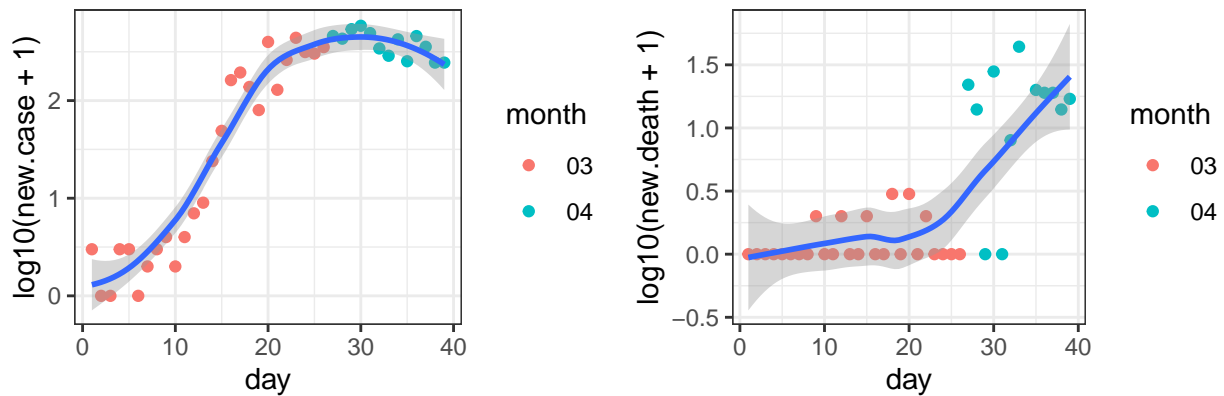
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-11

### Jefferson\_Louisiana



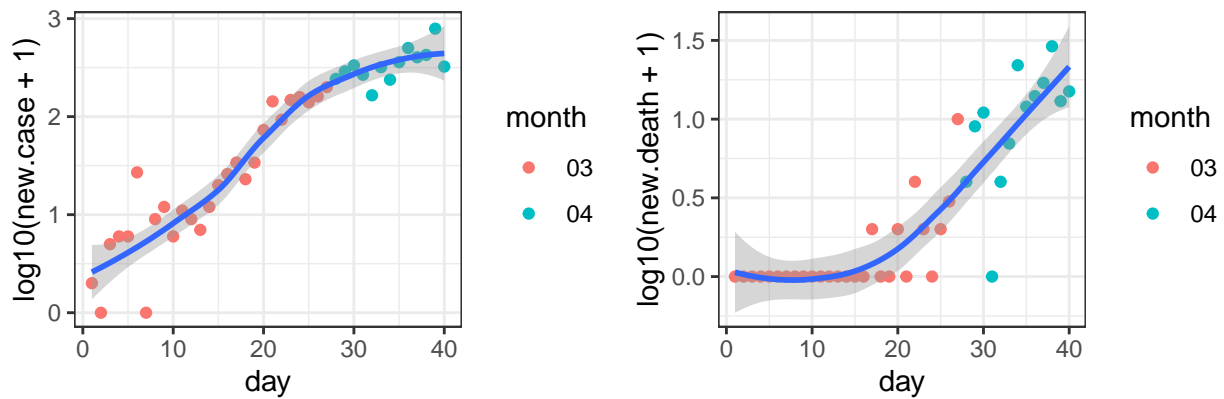
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-09

### Rockland\_New York



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

### Middlesex\_Massachusetts



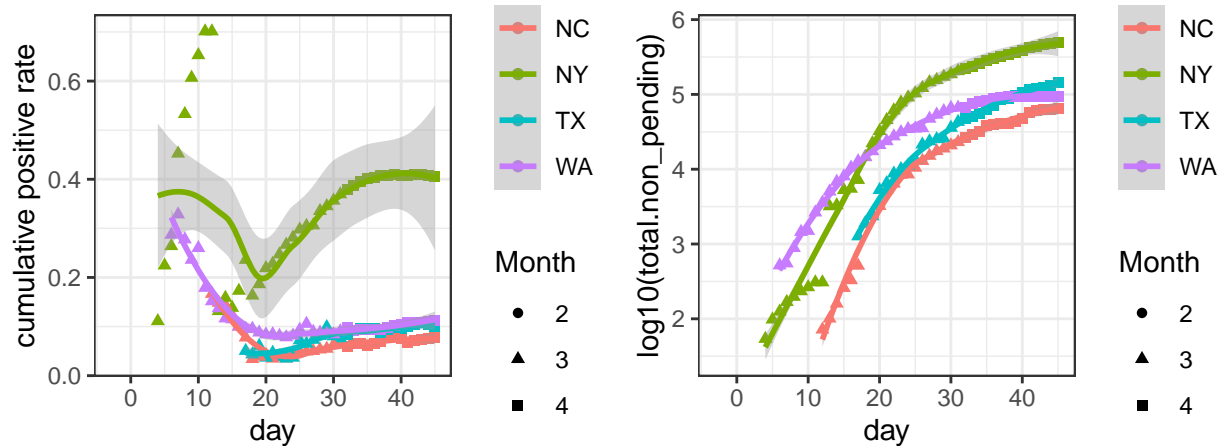
data source: <https://github.com/nytimes/covid-19-data>, day 1 is 03-05

## 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 0414: 0.11(WA) 0.10(TX) 0.41(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
```