

# ESTIMATING COVID-19'S $R_t$ IN REAL-TIME

THOMAS P. HARTE

## CONTENTS

1. Introduction	2
2. Bettencourt & Ribeiro's Approach	2
3. Choosing a Likelihood Function $Pr(k_t   R_t)$	3
4. Connecting $\lambda$ and $R_t$	5
5. Evaluating the Likelihood Function	5
6. Performing the Bayesian Update	6
7. Real-World Application to US Data	8
7.1. Setup	9
7.2. Running the Algorithm	10
7.2.1. Choosing the Gaussian $\sigma$ for $Pr(R_t   R_{t-1})$	10
7.3. Function for Calculating the Posteriors	11
7.4. The Result	11
7.5. Plotting in the Time Domain with Credible Intervals	12
7.6. Choosing the optimal $\sigma$	13
7.7. Compile Final Results	14
7.8. Plot All US States	16
7.9. Export Data to CSV	18
7.10. Standings	18
References	20
Appendix A. Build: Python packages	21

Original author: Kevin Systrom—April 17.

Source:

- (1) <https://rt.live/> (live page)
- (2) <https://github.com/k-sys/covid-19> (notebook);  
SHA1: cfba8cf5979da87adb7f57d6ff9ff3903d7a4f01

## 1. INTRODUCTION

In any epidemic,  $R_t$  is the measure known as the *effective reproduction number*. It's the number of people who become infected per infectious person at time  $t$ . The most well-known version of this number is the basic reproduction number:  $R_0$  when  $t = 0$ . However,  $R_0$  is a single measure that does not adapt with changes in behavior and restrictions.

As a pandemic evolves, increasing restrictions (or potential releasing of restrictions) changes  $R_t$ . Knowing the current  $R_t$  is essential. When  $R \gg 1$ , the pandemic will spread through a large part of the population. If  $R_t < 1$ , the pandemic will slow quickly before it has a chance to infect many people. The lower the  $R_t$ : the more manageable the situation. In general, any  $R_t < 1$  means things are under control.

The value of  $R_t$  helps us in two ways.

- (1) It helps us understand how effective our measures have been controlling an outbreak and
- (2) It gives us vital information about whether we should increase or reduce restrictions based on our competing goals of economic prosperity and human safety. [Well-respected epidemiologists argue](#) that tracking  $R_t$  is the only way to manage through this crisis.

Yet, today, we don't yet use  $R_t$  in this way. In fact, the only real-time measure I've seen has been for [Hong Kong](#). More importantly, it is not useful to understand  $R_t$  at a national level. Instead, to manage this crisis effectively, we need a local (state, county and/or city) granularity of  $R_t$ .

What follows is a solution to this problem at the US State level. It's a modified version of a solution created by [\[Bettencourt and Ribeiro, 2008\]](#) to estimate real-time  $R_t$  using a Bayesian approach. While this paper estimates a static  $R$  value, here we introduce a process model with Gaussian noise to estimate a time-varying  $R_t$ .

If you have questions, comments, or improvements feel free to get in touch: [hello@systrom.com](mailto:hello@systrom.com). And if it's not entirely clear, I'm not an epidemiologist. At the same time, data is data, and statistics are statistics and this is based on work by well-known epidemiologists so you can calibrate your beliefs as you wish. In the meantime, I hope you can learn something new as I did by reading through this example. Feel free to take this work and apply it elsewhere – internationally or to counties in the United States.

Additionally, a huge thanks to [Frank Dellaert](#) who suggested the addition of the process and to [Adam Lerer](#) who implemented the changes. Not only did I learn something new, it made the model much more responsive.

## 2. BETTENCOURT & RIBEIRO'S APPROACH

Every day, we learn how many more people have COVID-19. This new case count gives us a clue about the current value of  $R_t$ . We also, figure that the value of  $R_t$  today is related to the value of  $R_{t-1}$  (yesterday's value) and every previous value of  $R_{t-m}$  for that matter.

With these insights, the authors use [Bayes' rule](#) to update their beliefs about the true value of  $R_t$  based on how many new cases have been reported each day.

This is Bayes' Theorem as we'll use it:

$$\Pr(R_t | k) = \frac{\Pr(k | R_t) \cdot \Pr(R_t)}{\Pr(k)}$$

This says that, having seen  $k$  new cases, we believe the distribution of  $R_t$  is equal to:

- The **likelihood** of seeing  $k$  new cases given  $R_t$  times ...
- The **prior** beliefs of the value of  $\Pr(R_t)$  without the data ...
- divided by the probability of seeing this many cases in general.

This is for a single day. To make it iterative: every day that passes, we use yesterday's prior  $\Pr(R_{t-1})$  to estimate today's prior  $\Pr(R_t)$ . We will assume the distribution of  $R_t$  to be a Gaussian centered around  $R_{t-1}$ , so

$$\Pr(R_t | R_{t-1}) = \mathcal{N}(R_{t-1}, \sigma),$$

where  $\sigma$  is a hyperparameter (see below on how we estimate  $\sigma$ ). So on day one:

$$\Pr(R_1 | k_1) \propto \Pr(R_1) \cdot \mathcal{L}(R_1 | k_1)$$

On day two:

$$\begin{aligned} \Pr(R_2 | k_1, k_2) &\propto \Pr(R_2) \cdot \mathcal{L}(R_2 | k_2) \\ &= \sum_{R_1} \Pr(R_1 | k_1) \cdot \Pr(R_2 | R_1) \cdot \mathcal{L}(R_2 | k_2) \end{aligned}$$

etc.

### 3. CHOOSING A LIKELIHOOD FUNCTION $\Pr(k_t | R_t)$

A likelihood function says how likely we are to see  $k$  new cases, given a value of  $R_t$ .

Any time you need to model “arrivals” over some time period of time, statisticians like to use the [Poisson Distribution](#). Given an average arrival rate of  $\lambda$  new cases per day, the probability of seeing  $k$  new cases is distributed according to the Poisson distribution:

$$\Pr(k | \lambda) = \frac{\lambda^k e^{-\lambda}}{k!}$$

```

1 # Column vector of k
2 k = np.arange(0, 70)[:, None]
3
4 # Different values of Lambda
5 lambdas = [10, 20, 30, 40]
6
7 # Evaluated the Probability Mass Function (remember: poisson is discrete)
8 y = sps.poisson.pmf(k, lambdas)
9
10 # Show the resulting shape
11 print(y.shape)

```

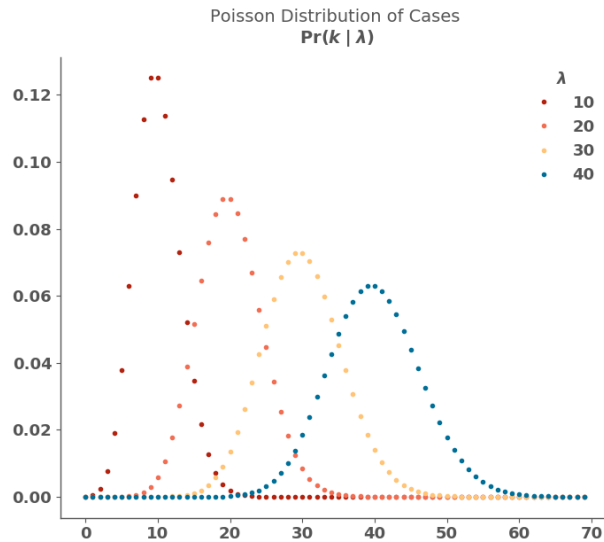
```
(70, 4)
```

**Note:** this was a terse expression which makes it tricky. All I did was to make  $k$  a column. By giving it a column for  $k$  and a “row” for lambda it will evaluate the pmf over both and produce an array that has  $k$  rows and lambda columns. This is an efficient way of producing many distributions all at once, and **you will see it used again below!**

```

1 def plot_it(filename):
2     fig, ax = plt.subplots(figsize=(8, 6.5))
3     ax.set(title='Poisson Distribution of Cases\n $\\Pr(k|\\lambda)$')
4     plt.plot(k, y,
5             marker='o',
6             markersize=3,
7             lw=0
8             )
9     plt.legend(title="$\\lambda$", labels=lambdas)
10    fig.savefig(filename)
11    return filename
12
13 plot_it(filename='./img/systrom-1.png')

```



The Poisson distribution says that if you think you're going to have  $\lambda$  cases per day, you'll probably get that many, plus or minus some variation based on chance.

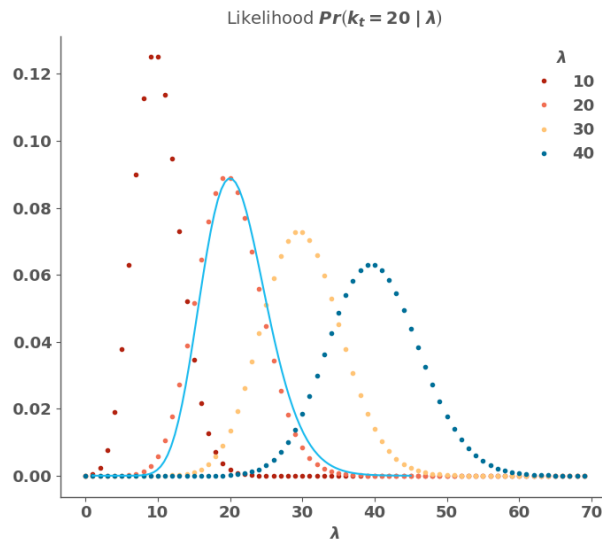
But in our case, we know there have been  $k$  cases and we need to know what value of  $\lambda$  is most likely. In order to do this, we fix  $k$  in place while varying  $\lambda$ . **This is called the likelihood function.**

For example, imagine we observe  $k = 20$  new cases, and we want to know how likely each  $\lambda$  is:

```

1 k = 20
2 lam = np.linspace(1, 45, 90)
3 likelihood = pd.Series(
4     data=sps.poisson.pmf(k, lam),
5     index=pd.Index(lam, name='lambda'),
6     name='likelihood'
7 )
8
9 def plot_it(filename):
10     likelihood.plot(
11         title=r'Likelihood  $\Pr(k=20 | \lambda)$ ',
12         figsize=(8, 6.5)
13     )
14     plt.legend(title="lambda", labels=lambdas)
15     plt.savefig(filename)
16     return filename
17
18 plot_it(filename='./img/systrom-2.png')

```



This says that if we see 20 cases, the most likely value of  $\lambda$  is (not surprisingly) 20. But we're not certain: it's possible  $\lambda$  was 21 or 17 and saw 20 new cases by chance alone. It also says that it's unlikely  $\lambda$  was 40 and we saw 20.

Great. We have  $Pr(\lambda_t | k_t)$  which is parameterized by  $\lambda$  but we were looking for  $Pr(k_t | R_t)$  which is parameterized by  $R_t$ . We need to know the relationship between  $\lambda$  and  $R_t$

#### 4. CONNECTING $\lambda$ AND $R_t$

**The key insight to making this work is to realize there's a connection between  $R_t$  and  $\lambda$ .** The derivation is beyond the scope of this notebook, but here it is:

$$\lambda = k_{t-1} e^{\gamma(R_t-1)}$$

where  $\gamma$  is the reciprocal of the serial interval (about 7 days for COVID19). Since we know every new case count on the previous day, we can now reformulate the likelihood function as a Poisson parameterized by fixing  $k$  and varying  $R_t$ .

$$Pr(k | R_t) = \frac{\lambda^k e^{-\lambda}}{k!}$$

#### 5. EVALUATING THE LIKELIHOOD FUNCTION

To continue our example, let's imagine a sample of new case counts  $k$ . What is the likelihood of different values of  $R_t$  on each of those days?

```

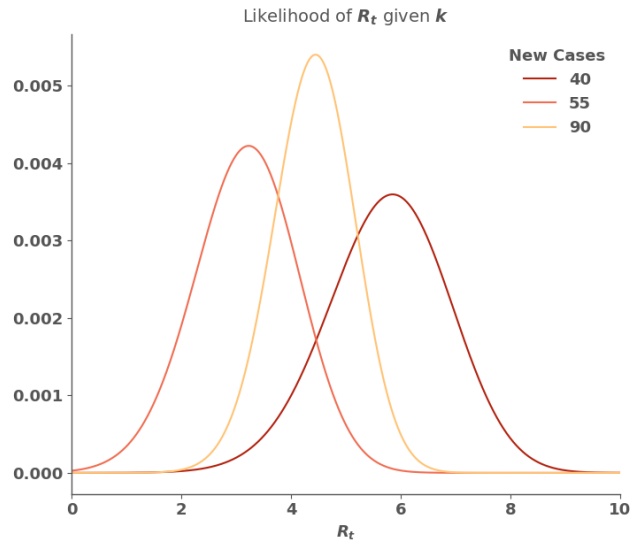
1 k = np.array([20, 40, 55, 90])
2
3 # We create an array for every possible value of Rt
4 R_T_MAX = 12
5 r_t_range = np.linspace(0, R_T_MAX, R_T_MAX*100+1)
6
7 # Gamma is 1/serial interval
8 # https://wwwnc.cdc.gov/eid/article/26/7/20-0282_article
9 # https://www.nejm.org/doi/full/10.1056/NEJMoa2001316
10 GAMMA = 1/7
11
12 # Map Rt into lambda so we can substitute it into the equation below
13 # Note that we have N-1 lambdas because on the first day of an outbreak
14 # you do not know what to expect.
15 lam = k[:-1] * np.exp(GAMMA * (r_t_range[:, None] - 1))

```

```

16
17 # Evaluate the likelihood on each day and normalize sum of each day to 1.0
18 likelihood_r_t = sps.poisson.pmf(k[1:], lam)
19 likelihood_r_t /= np.sum(likelihood_r_t, axis=0)
20
21 def plot_it(filename):
22     # Plot it
23     ax = pd.DataFrame(
24         data = likelihood_r_t,
25         index = r_t_range
26     ).plot(
27         title='Likelihood of $R_t$ given $k$',
28         xlim=(0,10),
29         figsize=(8,6.5)
30     )
31     ax.legend(labels=k[1:], title='New Cases')
32     ax.set_xlabel('$R_t$');
33     plt.savefig(filename)
34     return filename
35
36 plot_it(filename='./img/systrom-3.png')

```



You can see that each day we have a independent guesses for  $R_t$ . The goal is to combine the information we have about previous days with the current day. To do this, we use Bayes' theorem.

## 6. PERFORMING THE BAYESIAN UPDATE

To perform the Bayesian update, we need to multiply the likelihood by the prior (which is just the previous day's likelihood without our Gaussian update) to get the posteriors. Let's do that using the cumulative product of each successive day:

```

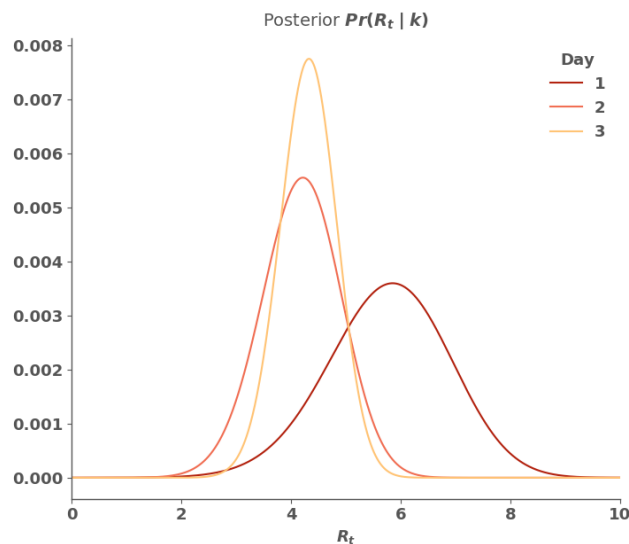
1 posteriors = likelihood_r_t.cumprod(axis=1)
2 posteriors = posteriors / np.sum(posteriors, axis=0)
3
4 columns = pd.Index(range(1, posteriors.shape[1]+1), name='Day')
5 posteriors = pd.DataFrame(
6     data = posteriors,
7     index = r_t_range,
8     columns = columns
9 )
10
11 def plot_it(filename):
12     ax = posteriors.plot(
13

```

```

14     title='Posterior  $\Pr(R_t | k)$ ',
15     xlim=(0,10),
16     figsize=(8,6.5)
17 )
18 ax.legend(title='Day')
19 ax.set_xlabel('$R_t$');
20 plt.savefig(filename)
21 return filename
22
23 plot_it(filename='./img/systrom-4.png')

```



Notice how on Day 1, our posterior matches Day 1's likelihood from above? That's because we have no information other than that day. However, when we update the prior using Day 2's information, you can see the curve has moved left, but not nearly as left as the likelihood for Day 2 from above. This is because Bayesian updating uses information from both days and effectively averages the two. Since Day 3's likelihood is in between the other two, you see a small shift to the right, but more importantly: a narrower distribution. We're becoming **more** confident in our beliefs of the true value of  $R_t$ .

From these posteriors, we can answer important questions such as :

“What is the most likely value of  $R_t$  each day?”

```

1 most_likely_values = posteriors.idxmax(axis=0)
2 print(most_likely_values)

```

```

Day
1    5.85
2    4.22
3    4.33
dtype: float64

```

We can also obtain the [highest density intervals](#) for  $R_t$ :

```

1 hdi = highest_density_interval(posteriors, debug=True)
2 print( f'\n{hdi.tail()}' )

```

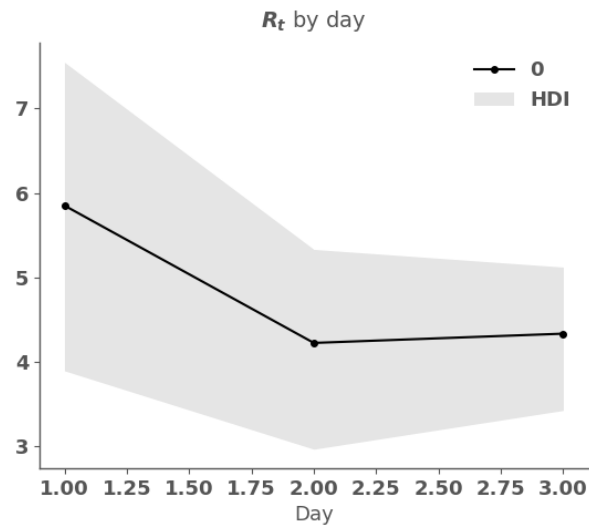
	Low_90	High_90
Day		
1	3.89	7.55
2	2.96	5.33
3	3.42	5.12

Finally, we can plot both the most likely values for  $R_t$  and the HDIs over time. This is the most useful representation as it shows how our beliefs change with every day.

```

1 def plot_it(filename):
2     # need a pd.DataFrame here
3     most_likely = pd.DataFrame(
4         most_likely_values
5     )
6     fig = plt.figure(figsize=(8, 6.5))
7     ax = most_likely.plot(
8         marker='o',
9         label='Most Likely',
10        title=f'$R_t$ by day',
11        c='k',
12        markersize=4,
13    )
14    ax.fill_between(
15        hdi.index,
16        hdi['Low_90'],
17        hdi['High_90'],
18        color='k',
19        alpha=.1,
20        lw=0,
21        label='HDI'
22    )
23    ax.legend();
24    plt.savefig(filename)
25    return filename
26
27 plot_it(filename='./img/systrom-5.png')

```



We can see that the most likely value of  $R_t$  changes with time and the highest-density interval narrows as we become more sure of the true value of  $R_t$  over time. Note that since we only had four days of history, I did not apply the process to this sample. Next, however, we'll turn to a real-world application where this process is necessary.

## 7. REAL-WORLD APPLICATION TO US DATA



### 7.1. Setup. Load US state case data from [CovidTracking.com](https://covidtracking.com)

```

1 url = 'https://covidtracking.com/api/v1/states/daily.csv'
2 states = pd.read_csv(
3     url,
4     usecols=['date', 'state', 'positive'],
5     parse_dates=['date'],
6     index_col=['state', 'date'],
7     squeeze=True
8 ).sort_index()

```

Taking a look at the state, we need to start the analysis when there are a consistent number of cases each day. Find the last zero new case day and start on the day after that.

Also, case reporting is very erratic based on testing backlogs, etc. To get the best view of the 'true' data we can, I've applied a gaussian filter to the time series. This is obviously an arbitrary choice, but you'd imagine the real world process is not nearly as stochastic as the actual reporting.

```

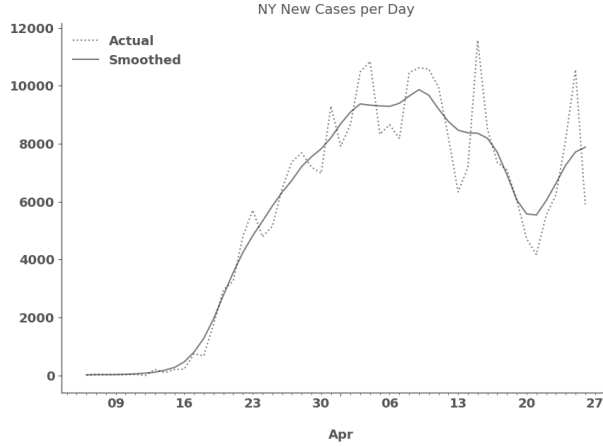
1 state_name = 'NY'
2
3 cases = states.xs(state_name).rename(f'{state_name} cases')
4 original, smoothed = prepare_cases(cases)

```

```

1 def make_plottable_df(series, name):
2     date = pd.to_datetime(
3         series.index, infer_datetime_format=True,
4         errors='coerce'
5     )
6     df = pd.DataFrame(
7         {name: series.values}
8     )
9     df.index = date
10    return df
11
12 def plot_it(filename):
13     df_o = make_plottable_df(original, 'Actual')
14     df_s = make_plottable_df(smoothed, 'Smoothed')
15     df = pd.concat([df_o, df_s], axis=1)
16     fig, ax = plt.subplots()
17     for cn in df.keys(): df[cn].plot(
18         ax=ax,
19         c='k',
20         linestyle = ':' if cn == 'Actual' else '-',
21         alpha=.5,
22         title=f'{state_name} New Cases per Day',
23         legend=True,
24         figsize=(10, 6.5)
25     )
26     ax.get_figure().set_facecolor('w')
27     plt.savefig(filename)
28     return filename
29
30 plot_it(filename='./img/systrom-6-a.png')

```



## 7.2. Running the Algorithm.

7.2.1. *Choosing the Gaussian  $\sigma$  for  $\Pr(R_t | R_{t-1})$ .* Note: you can safely skip this section if you trust that we chose the right value of  $\sigma$  for the process below. Otherwise, read on.

The original approach simply selects yesterday's posterior as today's prior. While intuitive, doing so doesn't allow for our belief that the value of  $R_t$  has likely changed from yesterday. To allow for that change, we apply Gaussian noise to the prior distribution with some standard deviation  $\sigma$ . The higher  $\sigma$  the more noise and the more we will expect the value of  $R_t$  to drift each day. Interestingly, applying noise on noise iteratively means that there will be a natural decay of distant posteriors. This approach has a similar effect of windowing, but is more robust and doesn't arbitrarily forget posteriors after a certain time like my previous approach. Specifically, windowing computed a fixed  $R_t$  at each time  $t$  that explained the surrounding  $w$  days of cases, while the new approach computes a series of  $R_t$  values that explains all the cases, assuming that  $R_t$  fluctuates by about  $\sigma$  each day.

However, there's still an arbitrary choice: what should  $\sigma$  be? Adam Lerer pointed out that we can use the process of maximum likelihood to inform our choice. Here's how it works:

Maximum likelihood says that we'd like to choose a  $\sigma$  that maximizes the likelihood of seeing our data  $k$ :  $P(k|\sigma)$ . Since  $\sigma$  is a fixed value, let's leave it out of the notation, so we're trying to maximize  $P(k)$  over all choices of  $\sigma$ .

Since  $P(k) = P(k_0, k_1, \dots, k_t) = P(k_0)P(k_1) \dots P(k_t)$  we need to define  $P(k_t)$ . It turns out this is the denominator of Bayes rule:

$$\Pr(R_t | k_t) = \frac{\Pr(k_t | R_t) \Pr(R_t)}{\Pr(k_t)}$$

To calculate it, we notice that the numerator is actually just the joint distribution of  $k$  and  $R$ :

$$\Pr(k_t, R_t) = \Pr(k_t | R_t) \Pr(R_t)$$

We can marginalize the distribution over  $R_t$  to get  $\Pr(k_t)$ :

$$\Pr(k_t) = \sum_{R_t} \Pr(k_t | R_t) \Pr(R_t)$$

So, if we sum the distribution of the numerator over all values of  $R_t$ , we get  $\Pr(k_t)$ . And since we're calculating that anyway as we're calculating the posterior, we'll just keep track of it separately.

Since we're looking for the value of  $\sigma$  that maximizes  $\Pr(k)$  overall, we actually want to maximize:

$$\prod_{t,i} p(k_{ti})$$

where  $t$  are all times and  $i$  is each state.

Since we're multiplying lots of tiny probabilities together, it can be easier (and less error-prone) to take the log of the values and add them together. Remember that  $\log ab = \log a + \log b$ . And since logarithms are monotonically increasing, maximizing the sum of the log of the probabilities is the same as maximizing the product of the non-logarithmic probabilities for any choice of  $\sigma$ .

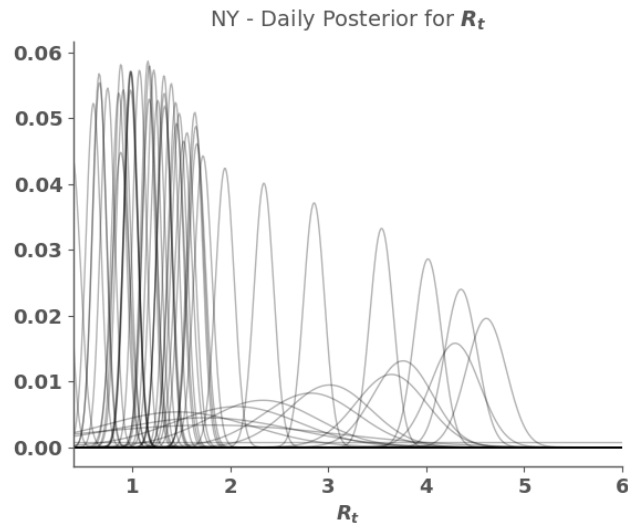
**7.3. Function for Calculating the Posteriors.** To calculate the posteriors we follow these steps:

- (1) Calculate  $\lambda$ , the expected arrival rate for every day's poisson process;
- (2) Calculate each day's likelihood distribution over all possible values of  $R_t$ ;
- (3) Calculate the process matrix based on the value of  $\sigma$  we discussed above;
- (4) Calculate our initial prior because our first day does not have a previous day from which to take the posterior
  - Based on [info from the cdc](#) we will choose a Gamma with mean 7.
- (5) Loop from day 1 to the end, doing the following:
  - Calculate the prior by applying the Gaussian to yesterday's prior.
  - Apply Bayes' rule by multiplying this prior and the likelihood we calculated in step 2.
  - Divide by the probability of the data (also Bayes' rule)

```
1 # Note that we're fixing sigma to a value just for the example
2 posteriors, log_likelihood = get_posteriors(smoothed, sigma=.25, gamma=GAMMA)
```

**7.4. The Result.** Below you can see every day (row) of the posterior distribution plotted simultaneously. The posteriors start without much confidence (wide) and become progressively more confident (narrower) about the true value of  $R_t$ .

```
1 def plot_it(filename):
2     ax = posteriors.plot(
3         title=f'{state_name} - Daily Posterior for $R_t$',
4         legend=False,
5         lw=1,
6         c='k',
7         alpha=.3,
8         xlim=(0.4,6)
9     )
10    ax.set_xlabel('$R_t$')
11    plt.savefig(filename)
12    return filename
13
14 plot_it(filename='./img/systrom-7.png')
```



**7.5. Plotting in the Time Domain with Credible Intervals.** Since our results include uncertainty, we'd like to be able to view the most likely value of  $R_t$  along with its highest-density interval.

```

1 hdis = highest_density_interval(posterior, p=.9)
2 most_likely = posterior.idxmax().rename('ML')
3
4 # Look into why you shift -1
5 result = pd.concat([most_likely, hdis], axis=1)
6 print( f'\n{result.tail()}' )

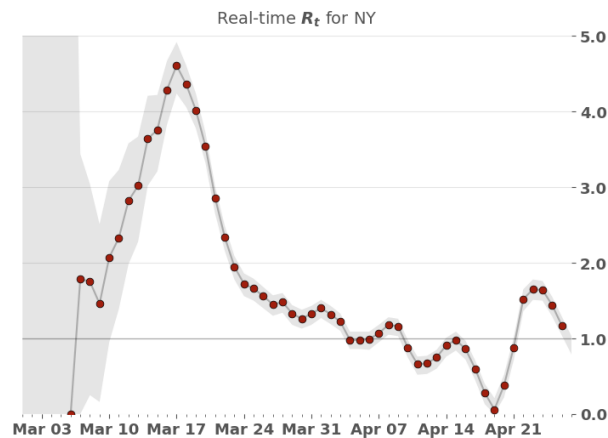
```

	ML	Low_90	High_90
date			
2020-04-22	1.52	1.36	1.65
2020-04-23	1.65	1.51	1.78
2020-04-24	1.64	1.50	1.76
2020-04-25	1.44	1.29	1.55
2020-04-26	1.17	1.04	1.29

```

1 def plot_it(filename):
2     fig, ax = plt.subplots(figsize=(600/72,400/72))
3     plot_rt(result, ax, state_name, fig)
4     ax.set_title(f'Real-time $R_t$ for {state_name}')
5     ax.xaxis.set_major_locator(mdates.WeekdayLocator())
6     ax.xaxis.set_major_formatter(mdates.DateFormatter('%b %d'))
7     plt.savefig(filename)
8     return filename
9
10 plot_it(filename='./img/systrom-8.png')

```



**7.6. Choosing the optimal  $\sigma$ .** In the previous section we described choosing an optimal  $\sigma$ , but we just assumed a value. But now that we can evaluate each state with any sigma, we have the tools for choosing the optimal  $\sigma$ .

Above we said we'd choose the value of  $\sigma$  that maximizes the likelihood of the data  $P(k)$ . Since we don't want to overfit on any one state, we choose the sigma that maximizes  $P(k)$  over every state. To do this, we add up all the log likelihoods per state for each value of sigma then choose the maximum.

Note: this takes a while!

```
1 states.loc[('MA', pd.Timestamp('2020-04-20'))] = 39643
```

```
1 print( f"\n{states.xs('MA').diff().tail()}" )
```

```
.
date
2020-04-22    1745.0
2020-04-23    3079.0
2020-04-24         0.0
2020-04-25    7325.0
2020-04-26    1590.0
Name: positive, dtype: float64
```

```
1 sigmas = np.linspace(1/20, 1, 20)
2
3 targets = ~states.index.get_level_values('state').isin(FILTERED_REGION_CODES)
4 states_to_process = states.loc[targets]
5
6 if not 'results' in locals() or not len(results):
7     results = {}
8     for state_name, cases in states_to_process.groupby(level='state'):
9         print(state_name)
10        new, smoothed = prepare_cases(cases, cutoff=25)
11        if len(smoothed) == 0:
12            new, smoothed = prepare_cases(cases, cutoff=10)
13        result = {}
14        # Holds all posteriors with every given value of sigma
15        result['posteriors'] = []
16        # Holds the log likelihood across all k for each value of sigma
17        result['log_likelihoods'] = []
18        for sigma in sigmas:
19            posteriors, log_likelihood = get_posteriors(smoothed, sigma=sigma)
20            result['posteriors'].append(posteriors)
21            result['log_likelihoods'].append(log_likelihood)
```

```

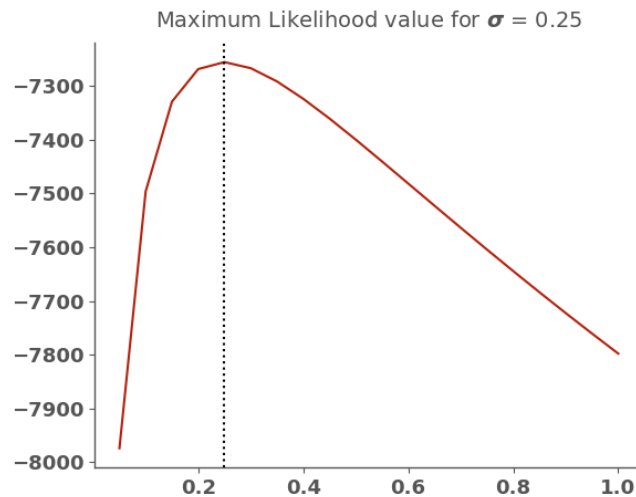
22     # Store all results keyed off of state name
23     results[state_name] = result
24
25     print('Done.')
```

```
Done.
```

Now that we have all the log likelihoods, we can sum for each value of sigma across states, graph it, then choose the maximum.

```

1  # Each index of this array holds the total of the log likelihoods for
2  # the corresponding index of the sigmas array.
3  total_log_likelihoods = np.zeros_like(sigmas)
4
5  # Loop through each state's results and add the log likelihoods to the running total.
6  for state_name, result in results.items():
7      total_log_likelihoods += result['log_likelihoods']
8
9  # Select the index with the largest log likelihood total
10 max_likelihood_index = total_log_likelihoods.argmax()
11
12 # Select the value that has the highest log likelihood
13 sigma = sigmas[max_likelihood_index]
14 def plot_it(filename):
15     fig, ax = plt.subplots()
16     ax.set_title(f"Maximum Likelihood value for  $\sigma$  = {sigma:.2f}");
17     ax.plot(sigmas, total_log_likelihoods)
18     ax.axvline(sigma, color='k', linestyle=":")
19     plt.savefig(filename)
20     return filename
21
22 plot_it(filename='./img/systrom-9.png')
```



**7.7. Compile Final Results.** Given that we've selected the optimal  $\sigma$ , let's grab the precalculated posterior corresponding to that value of  $\sigma$  for each state. Let's also calculate the 90% and 50% highest density intervals (this takes a little while) and also the most likely value.

```

1  if not 'final_results' in locals() or final_results is None:
2      final_results = None
3
4      for state_name, result in results.items():
5          print(f'{state_name} ')
6          posteriors = result['posteriors'][max_likelihood_index]
```

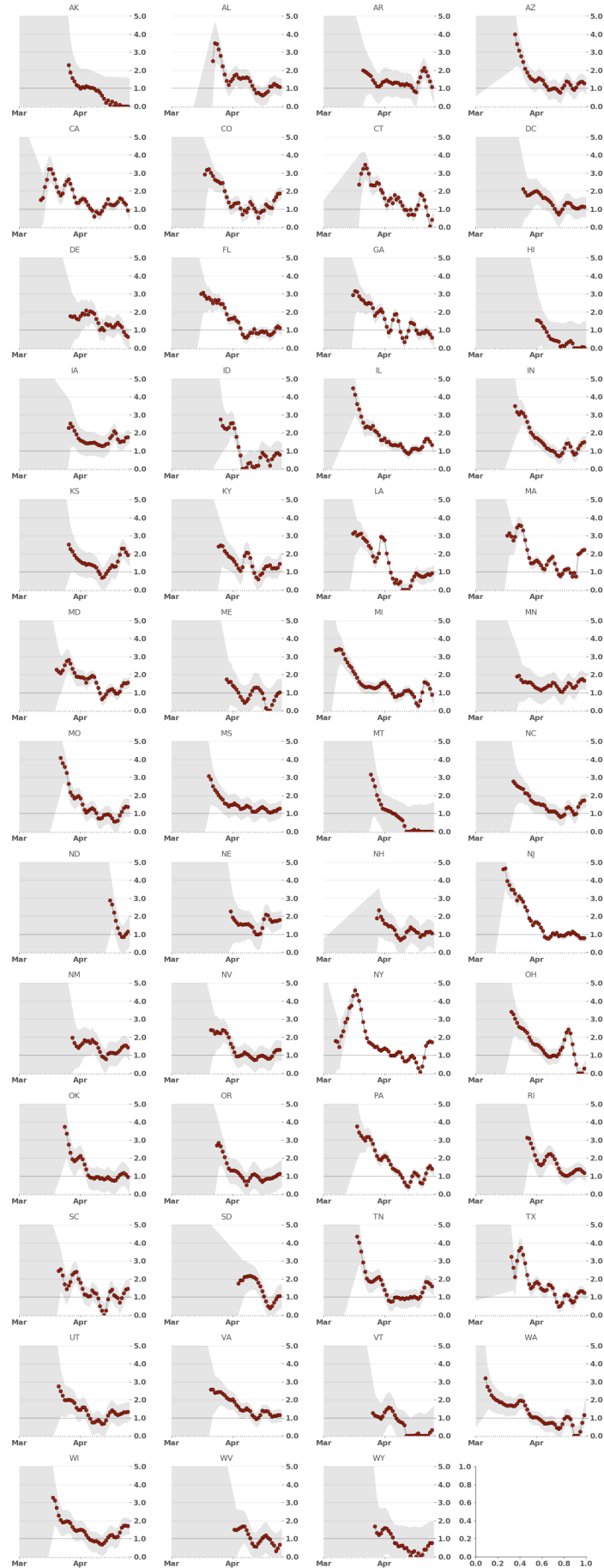
```
7     hdis_90 = highest_density_interval(posterior, p=.9)
8     hdis_50 = highest_density_interval(posterior, p=.5)
9     most_likely = posterior.idxmax().rename('ML')
10    result = pd.concat([most_likely, hdis_90, hdis_50], axis=1)
11    if final_results is None:
12        final_results = result
13    else:
14        final_results = pd.concat([final_results, result])
15
16    print('Done.')
```

Done.

```
1 ncols = 4
2 nrows = int(np.ceil(len(results) / ncols))
3
4 def plot_it(filename):
5     fig, axes = plt.subplots(nrows=nrows, ncols=ncols, figsize=(15, nrows*3))
6     for i, (state_name, result) in enumerate(final_results.groupby('state')): plot_rt(result.iloc[1:], axes.flat[i], state_name, fig)
7     fig.tight_layout()
8     fig.set_facecolor('w')
9     plt.savefig(filename)
10    return filename
11
12 plot_it(filename='./img/systrom-10.png')
```

## 7.8. Plot All US States.





```

1 # Since we now use a uniform prior, the first datapoint is pretty bogus, so just truncating it here
2 final_results = final_results.groupby('state').apply(lambda x: x.iloc[1:].droplevel(0))

```

## 7.9. Export Data to CSV. Create a data folder:

```

1 mkdir -p ./data

```

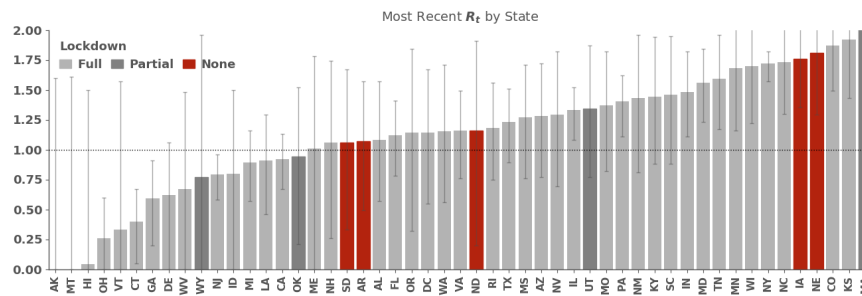
Export the data the data folder:

```

1 filtered = final_results.index.get_level_values(0).isin(FILTERED_REGION_CODES)
2 mr = final_results.loc[~filtered].groupby(level=0)[['ML', 'High_90', 'Low_90']].last()
3
4 def plot_it(filename):
5     mr.sort_values('ML', inplace=True)
6     plot_standings(mr, figsize=(16, 4.5));
7     plt.savefig(filename)
8     return filename
9
10 plot_it(filename='./img/systrom-11.png')

```

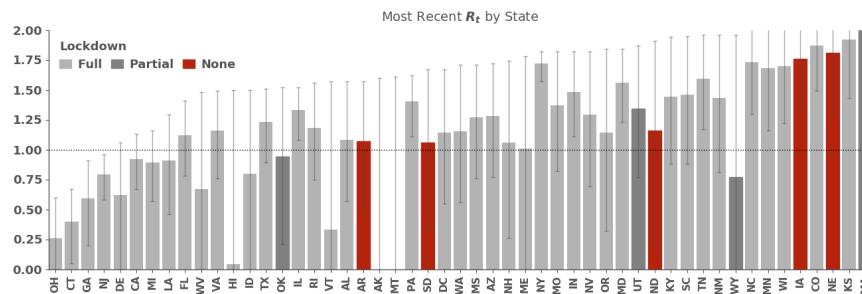
## 7.10. Standings.



```

1 def plot_it(filename):
2     mr.sort_values('High_90', inplace=True)
3     plot_standings(mr, figsize=(16, 4.5));
4     plt.savefig(filename)
5     return filename
6
7 plot_it(filename='./img/systrom-12.png')

```



```

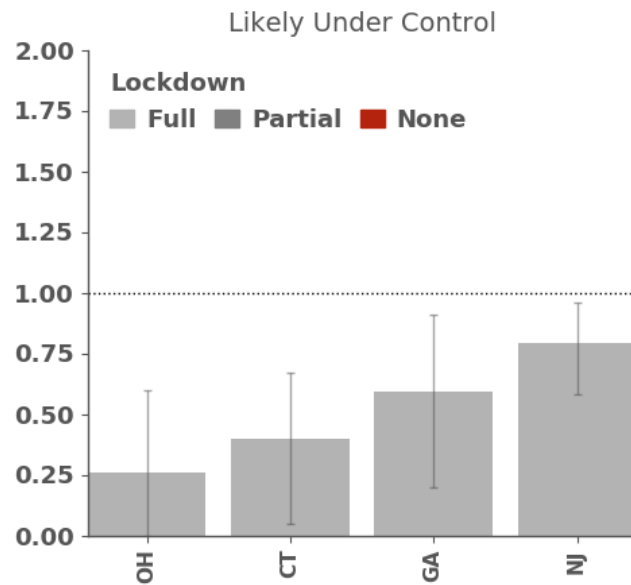
1 def plot_it(filename):
2     show = mr[mr.High_90.le(1)].sort_values('ML')
3     fig, ax = plot_standings(show, figsize=(5.25, 4.5), title='Likely Under Control');
4     plt.savefig(filename)

```

```

5     return filename
6
7     plot_it(filename='./img/systrom-13.png')

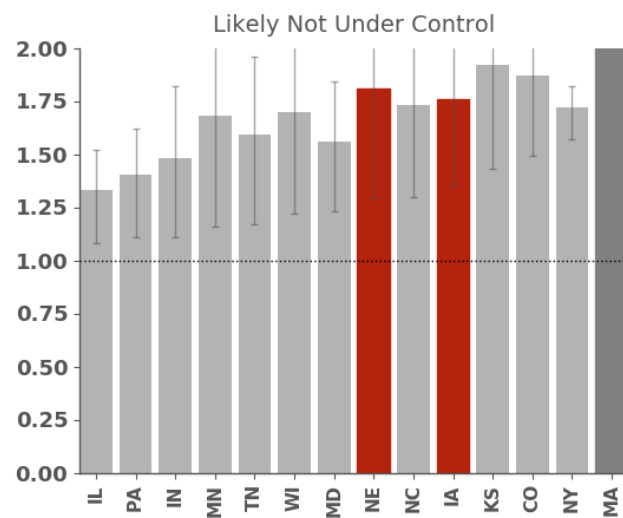
```



```

1 show = mr[mr.High_90.le(1)].sort_values('ML')
2 fig, ax = plot_standings(show, title='Likely Under Control');
3
4 def plot_it(filename):
5     show = mr[mr.Low_90.ge(1.0)].sort_values('Low_90')
6     fig, ax = plot_standings(show, figsize=(6, 4.5), title='Likely Not Under Control');
7     ax.get_legend().remove()
8     plt.savefig(filename)
9     return filename
10
11 plot_it(filename='./img/systrom-14.png')

```



## REFERENCES

[Bettencourt and Ribeiro, 2008] Bettencourt, L. M. A. and Ribeiro, R. M. (2008). Real Time Bayesian Estimation of the Epidemic Potential of Emerging Infectious Diseases. *PLOS ONE* 5.

## APPENDIX A. BUILD: PYTHON PACKAGES

```
pd.show_versions()
```

**INSTALLED VERSIONS**

-----

commit: None  
python: 3.7.1.final.0  
python-bits: 64  
OS: Linux  
OS-release: 5.3.0-46-generic  
machine: x86\_64  
processor: x86\_64  
byteorder: little  
LC\_ALL: None  
LANG: en\_US.UTF-8  
LOCALE: en\_US.UTF-8

pandas: 0.23.4  
pytest: 4.0.2  
pip: 18.1  
setuptools: 40.6.3  
Cython: 0.29.2  
numpy: 1.15.4  
scipy: 1.1.0  
pyarrow: None  
xarray: None  
IPython: 7.2.0  
sphinx: 1.8.2  
patsy: 0.5.1  
dateutil: 2.7.5  
pytz: 2018.7  
blosc: None  
bottleneck: 1.2.1  
tables: 3.4.4  
numexpr: 2.6.8  
feather: None  
matplotlib: 3.0.2  
openpyxl: 2.5.12  
xlrd: 1.2.0  
xlwt: 1.3.0  
xlsxwriter: 1.1.2  
lxml: 4.2.5  
bs4: 4.6.3  
html5lib: 1.0.1  
sqlalchemy: 1.2.15  
pymysql: None  
psycopg2: None  
jinja2: 2.10  
s3fs: None  
fastparquet: None  
pandas\_gbq: None  
pandas\_datareader: None