

# Tracking Reproductivity of COVID-19 Epidemic in China with Varying Coefficient SIR Model

HAOXUAN SUN<sup>1</sup>, YUMOU QIU<sup>2</sup>, HAN YAN<sup>3</sup>, YAXUAN HUANG<sup>4</sup>, YURU ZHU<sup>5</sup>, JIA GU<sup>5</sup>, AND  
SONG XI CHEN<sup>\*5,6</sup>

<sup>1</sup>Center for Data Science, Peking University, Beijing, China

<sup>2</sup>Department of Statistics, Iowa State University, Ames, Iowa, USA

<sup>3</sup>School of Mathematical Sciences, Sichuan University, Chengdu, Sichuan, China

<sup>4</sup>Yuanpei College, Peking University, Beijing, China

<sup>5</sup>Center for Statistical Science, Peking University, Beijing, China

<sup>6</sup>Guanghua School of Management, Peking University, Beijing China

## Abstract

We propose a varying coefficient Susceptible-Infected-Removal (vSIR) model that allows changing infection and removal rates for the latest corona virus (COVID-19) outbreak in China. The vSIR model together with proposed estimation procedures allow one to track the reproductivity of the COVID-19 through time and to assess the effectiveness of the control measures implemented since Jan 23 2020 when the city of Wuhan was lockdown followed by an extremely high level of self-isolation in the population. Our study finds that the reproductivity of COVID-19 had been significantly slowed down in the three weeks from January 27th to February 17th with 96.3% and 95.1% reductions in the effective reproduction numbers  $R$  among the 30 provinces and 15 Hubei cities, respectively. Predictions to the ending times and the total numbers of infected are made under three scenarios of the removal rates. The paper provides a timely model and associated estimation and prediction methods which may be applied in other countries to track, assess and predict the epidemic of the COVID-19 or other infectious diseases.

**Keywords** *epidemic assessment; estimation of basic reproductive number*

## 1 Introduction

The Corona Virus Disease 2019 (COVID-19) has created a profound public health emergency in China and has spread to 25 countries so far ([World Health Organization, 2020](#)). It has become an epidemic with more than 76,000 confirmed infections and 2,244 reported deaths worldwide as on February 20 2020. The COVID-19 is caused by a new corona viruses that is genetically similar to the viruses causing severe acute respiratory syndrome (SARS) and Middle East respiratory syndrome (MERS). Despite a relatively lower fatality rate comparing to SARS and MERS, the COVID-19 spreads faster and infects much more people than the SARS-03 outbreak.

The city of Wuhan, the origin of the outbreak, has been locked up to reduce population movement since January 23 in an effort to stop the spread of the epidemic, followed by more than 50 prefecture level cities (as on 8th of February) and countless number of towns and villages in China. A high percentage of the population are exercising self-isolation in their homes. The spring festival holiday period had been extended with all schools and universities closed and all students staying where they are indefinitely. The country is virtually in a stand-still, and the economy and people's livelihood have been severely affected by the epidemic.

---

\*The corresponding authors are Song Xi Chen (csx@gsm.pku.edu.cn) and Yumou Qiu (yumouqiu@iastate.edu).

There is an urgent need to assess the speed of the disease transmission and to check if the existing containment measures have successfully slowed down the spread of the disease or not. The Susceptible-Infected-Removal (SIR) model (Kermack and McKendrick, 1927) and its generalizations, for instance the Susceptible-Exposed-Infected-Removal (SEIR) model (Hethcote, 2000) with four or more compartments are commonly used to model the dynamics of infectious disease outbreaks. See (Becker, 1977; Becker and Britton, 1999; Yip and Chen, 1998; Ball and Clancy, 1993) for statistical estimation and inference for stochastic versions of the SIR model. SEIR models have been used to produce early results on COVID-19 in (Wu et al., 2020; Read et al., 2020; Tang et al., 2020), which produced the first three estimates of the basic reproduction number  $R_0$ : 2.68 by (Wu et al., 2020), 3.81 by (Read et al., 2020) and 6.47 by (Tang et al., 2020). The  $R_0$  is the expected number of infections by one infectious person over his/her infectious period at the start of the epidemic, which is closely connected to the effective reproduction number  $R_t$ . The latter  $R_t$  is the expected number of infections by one infected over infectious period at time  $t$  of the epidemic. Both  $R_0$  and  $R_t$  are key measures of an epidemic. For fixed coefficient models, if  $R_0 < 1$ , the epidemic will die down eventually with the speed of the decline depends on the size of  $R_0$ ; otherwise, the epidemic will explode until it runs out of its course.

The SEIR models that was employed in the above three cited works for the COVID-19 assume constant model coefficients, implying a constant regime of transmission during the course of the epidemic. This is idealistic for modeling COVID-19 as it cannot reflect the intervention measures by the authorities and the citizens, which should have made the infectious rate ( $\beta$ ) and the effective reproduction number ( $R_t$ ) varying with respect to time. Here, the effective reproductive number  $R_t$  is the average number of secondary infections made by each infectious case during an epidemic, which contrasts the basic reproductive number  $R_0$  that measures the average number of secondary infections at the beginning of an epidemic.

To reflect the changing dynamic regimes due to the strong government intervention and the self protective reactions by citizens, we propose a varying coefficient SIR (vSIR) model. The vSIR model is easy to be implemented via the locally weighted regression approach (Cleveland and Devlin, 1988) that produces estimates with desired smoothness, and yet is able to capture the changing dynamics of COVID-19's reproduction, with guaranteed statistical consistency and needed standard errors. The consistent estimator and its confidence interval are proposed for estimating the trend of  $R$ , assessing the effectiveness of infection control, and predicting the ending time and the final number of infection cases with 95% prediction intervals.

As COVID-19 is quickly spreading outside China, the vSIR model and the associated estimation and prediction methods may be applied to other countries to track, assess and predict the epidemic of the COVID-19 or other infectious diseases.

## 2 Main Results

By applying the vSIR model, we produce daily estimates of the infectious rate  $\beta(t)$  and the effective reproduction number  $R_t^D$  ( $t$  denotes time) based on three values of infectious duration  $D$ : 7, 10.5 and 14 days for 30 provinces and 15 major cities (including Wuhan) in Hubei province from January 21 or a later date between January 24-29 depending on the first confirmed case to February 17.

- Despite the total number of confirmed cases and the death are increasing, the spread of COVID-19 has shown a great slowing down in China within the two weeks from January 27 to February 17 as shown by 96.3% and 95.1% reductions in the effective reproduction number

$R_t$  among the 30 provinces and the 15 cities in Hubei, respectively.

- The average  $R_t^{14}$  (based on 14-day infectious duration) on January 27th was 6.14 (1.49) and 7.59 (2.38), respectively, for the 27 provinces and the 7 Hubei cities with confirmed cases by January 23rd. The numbers in the parentheses are the standard error. One week later on February 3rd, the  $R_t^{14}$  was averaged at 2.18 (0.67) for the 30 provinces and 2.84 (0.59) for the 15 Hubei cities, representing 64.5% and 62.6% reductions, respectively, over the 7 days. On February 10th, the average  $R_t^{14}$  dropped further to 0.86 (0.38) for the 30 provinces and 1.23 (0.55) for the 15 Hubei cities, which were either below or close to the critical threshold level 1.
- On February 17th, the average  $R_t^{14}$  has reached 0.23(0.15) and 0.37(0.24) for the 30 provinces and the 15 Hubei cities, with 22 provinces' and 8 Hubei cities'  $R_t^{14}$  being statistically significantly below 1 for more than 7 days. These indicate a further slowing down in the re-productivity of COVID-19 in China in the week from February 10 to 17.
- The profound slowing down in the reproductivity of COVID-19 can be attributed to a series of containment measures by the government and the public, which include cutting off Wuhan and other cities from January 23, a rapid public awareness of the epidemic and the extensive self protection taken and high level of self isolation at home exercised over a much extended Spring Festival holiday period.
- There are increasing numbers of provinces and cities in Hubei whose 14-day  $R_t$  has been statistically below 1, as detailed in Table 1, which would foreshadow the coming of the turning point for containment of the epidemic, if the control measures implemented since January 23 can be continued.
- If the recovery rate can be increased to 0.1 meaning the average recover time is 10 days after diagnosis, the number of infected patients  $I(t)$  will be dramatically reduced in March, and the epidemic will end in April for non-Hubei provinces and end in June for Hubei.

Table 1: The reproduction number  $R_t^D$  at two infectious durations: 10.5 and 14 days, for the 30 mainland provinces and 15 cities in Hubei province on February 10th with extended results on February 17th. The symbols + (−) indicate that the  $R_t^{14}$  was significantly above (below) 1 at 5% level of statistical significance, and the numbers inside the square brackets were the consecutive days the  $R_t^{14}$  were significantly below 1. The column  $\Delta R$  gives the percentages of decline in the  $R_t^{14}$  from the beginning of analysis to February 10th (the first two weeks of the analysis). The columns  $\Delta R(1^{st})$ ,  $\Delta R(2^{nd})$  and  $\Delta R(3^{rd})$  are the percentages of decline in the first week (January 27 to February 3rd), the second week (February 3-10), and the third week (February 10-17), respectively.

Province/City	$R_t^{10.5}$	$R_t^{14}$	$\Delta R$	$\Delta R(1^{st})$	$\Delta R(2^{nd})$	$\Delta R(3^{rd})$
Wuhan	1.99+	2.66+	58.7%	45.9%	23.7%	72.5%
Ezhou	1.64+	2.18+	80%	79.3%	3.6%	67.7%
Hubei	1.48+	1.98+	74.2%	58.2%	38.3%	69.7%
Tianmen	1.33+	1.78+	75%	67.4%	23.4%	52.5%
Guizhou	1.25+	1.67+	62.4%	9.3%	58.5%	91.5%
Xiantao	0.99	1.32	76.9%	46.4%	57%	71.9%
Heilongjiang	0.95	1.27+	81.8%	54.3%	60.3%	62.6%
Hebei	0.94	1.25+	85.4%	82.4%	16.7%	70.5%

(Continued on the next page)

Table 1 — Continued from the previous page

Province/City	$R_t^{10.5}$	$R_t^{14}$	$\Delta R$	$\Delta R(1^{st})$	$\Delta R(2^{nd})$	$\Delta R(3^{rd})$
Xinjiang	0.9	1.2+	75.6%	60.7%	37.9%	53.1%
Enshizhou	0.86–[5]	1.14+	74.1%	60.3%	34.6%	76%
Jingzhou	0.85–[1]	1.14+	84.8%	50.5%	69.3%	76.9%
Gansu	0.8	1.07	75.1%	47.8%	52.3%	100%
Jingmen	0.79–[1]	1.05	88.3%	77.1%	49.2%	84.1%
Huangshi	0.79–[1]	1.05	78.3%	31.2%	68.4%	83.6%
Anhui	0.74–[1]	0.99	88.3%	71.7%	58.7%	77.6%
Shanxi	0.74–[2]	0.98	86.7%	69.6%	56.1%	87.1%
Ningxia	0.73	0.97	84.9%	75.8%	37.6%	75.7%
Shandong	0.73–[2]	0.97	90.3%	84.5%	37.1%	80.3%
Jiangsu	0.72–[3]	0.96	87.1%	70.5%	56.1%	72.1%
Xianning	0.71–[4]	0.95	70.7%	21.6%	62.6%	33.6%
Shiyan	0.71–[1]	0.94	89.1%	72.4%	60.3%	67.5%
Jilin	0.7	0.93	80.4%	17.6%	76.1%	82.5%
Yichang	0.69–[3]	0.92	87%	56.6%	70%	75.7%
Huanggang	0.69–[3]	0.92	88.9%	60.7%	71.8%	89%
Tianjin	0.69–[4]	0.91	82.9%	51.8%	64.6%	64%
Hainan	0.68–[1]	0.91	80.9%	65.2%	45.2%	96.1%
Guangxi	0.66–[5]	0.88	81.6%	64.9%	47.8%	72.1%
Xiangyang	0.63–[3]	0.84–[1]	88.4%	57.9%	72.4%	80.8%
Sichuan	0.62–[5]	0.83–[2]	89.4%	78.5%	50.7%	47.6%
Jiangxi	0.61–[2]	0.82	90.8%	70.9%	68.5%	79.6%
Xiaogan	0.6–[1]	0.81	88.9%	61.5%	71.3%	50.1%
Hunan	0.57–[3]	0.76–[2]	91.5%	77%	63.2%	90.4%
Henan	0.56–[2]	0.75–[1]	93.2%	78.8%	67.8%	64.2%
Suizhou	0.52–[2]	0.69–[2]	88.2%	40.9%	80%	65.5%
Chongqing	0.51–[4]	0.68–[3]	90.4%	75%	61.6%	73.9%
Shaanxi	0.51–[3]	0.68–[2]	86.5%	63%	63.6%	72.9%
Neimenggu	0.49–[3]	0.66–[2]	82.4%	43%	69.1%	27.9%
Fujian	0.49–[6]	0.66–[4]	90.5%	76.2%	60%	76.9%
Guangdong	0.45–[3]	0.61–[2]	88.2%	54.4%	74.2%	62%
Liaoning	0.45–[6]	0.6–[2]	89.3%	72.7%	61%	81.7%
Beijing	0.45–[4]	0.6–[2]	90.3%	59.2%	76.2%	65.1%
Shanghai	0.34–[4]	0.46–[2]	92.1%	68%	75.2%	53.4%
Zhejiang	0.31–[4]	0.42–[3]	94.3%	77.9%	74.2%	86.3%
Yunnan	0.28–[7]	0.38–[5]	96.2%	86.8%	71.5%	30.8%
Qinghai	0.02–[4]	0.03–[3]	98.9%	-1.6%	98.9%	100%
Average (sd)	0.74 (0.35)	0.98 (0.47)	85.3%	64.2%	59%	71.6%

### 3 Time-varying coefficient SIR model

Let  $S(t)$ ,  $I(t)$  and  $R(t)$  be the counts of susceptible, infected and removed (including dead) persons in a given city or province at time  $t$ , respectively. Let  $N$  be the total population of

the city/province. We propose a varying coefficient Susceptible-Infected-Removed (vSIR) model for the conditional means of the Poisson increments  $\Delta I(t)$  and  $\Delta R(t)$  given  $I(t)$  and  $R(t)$ . This Poisson-vSIR framework permits estimating the parameters and the effective reproduction number  $R_t$  for the dynamics of COVID-19, which are then used for predicting the future spread of the disease.

The SIR model (Kermack and McKendrick, 1927) is a commonly used epidemiology model for the dynamic of susceptible  $S(t)$ , infected  $I(t)$  and recovered  $R(t)$  as a system of ordinary differential equations (ODEs). Here we consider a generalized version of the SIR model in that the infectious rate  $\beta$  and the removal rate  $\gamma$  may vary with respect to time so that the deterministic ODEs are

$$\begin{aligned}\frac{dS(t)}{dt} &= -\beta(t)I(t)\frac{S(t)}{N}, \\ \frac{dI(t)}{dt} &= \beta(t)I(t)\frac{S(t)}{N} - \gamma(t)I(t), \\ \frac{dR(t)}{dt} &= \gamma(t)I(t),\end{aligned}\tag{1}$$

where  $\beta(t)$  and  $\gamma(t)$  are unknown infection and the removal rate functions, respectively. Once an individual is removed, including dead, the individual can not return to the susceptible group.

The rationale for using a time-varying  $\beta(t)$  function, rather than a constant  $\beta$ , is that  $\beta(t)$  is the average rate of contact per unit time multiplied by the probability of disease transmission per contact between a susceptible and an infectious subject. Due to an increasing public awareness of the epidemic and the control measures as mentioned earlier, both the transmission probability and the contact rate have been largely reduced. These favor for a time-varying  $\beta(t)$ , which are also confirmed by the sharp declined in  $R_t^D = \beta(t)D$ , where  $D$  denote the infectious durations in Figures 1 and 2. The removal rate also changes over time as treatments improve over time as shown in Figure 3. However, Figure 3 shows  $\gamma(t)$  is much slowly changing for most of the provinces, which led us to treat  $\gamma(t) = \gamma$  at the early stage of the outbreak, whose value gradually increased as the recover rate improved as the time progress and better treatments are available.

The deterministic vSIR model as specified by the ODEs in (1) specifies the conditional means of the Poisson increments  $\Delta I(t)$  and  $\Delta R(t)$  given  $S(t)$ ,  $I(t)$  and  $R(t)$  at each discrete time point  $t$ . This conditional mean specification leads to a Poisson-vSIR model framework, which can be used to construct conditional likelihood for  $(\beta(t), \gamma(t))$  over moving time windows and leads to statistical inference for the effective reproduction number estimation and its standard error. The Poisson-vSIR framework is also the basis for the bootstrap re-sampling algorithm that we will propose for generating predictive intervals.

SEIR model is an extension of SIR with an added compartment  $E$  for the exposed between  $S$  and  $I$ . A time-varying SEIR model (vSEIR) satisfies the ODEs

$$\begin{aligned}\frac{dS(t)}{dt} &= -\beta(t)I(t)s(t), \\ \frac{dE(t)}{dt} &= \beta(t)I(t)s(t) - \alpha(t)E(t), \\ \frac{dI(t)}{dt} &= \alpha(t)E(t) - \gamma(t)I(t), \\ \frac{dR(t)}{dt} &= \gamma(t)I(t),\end{aligned}\tag{2}$$

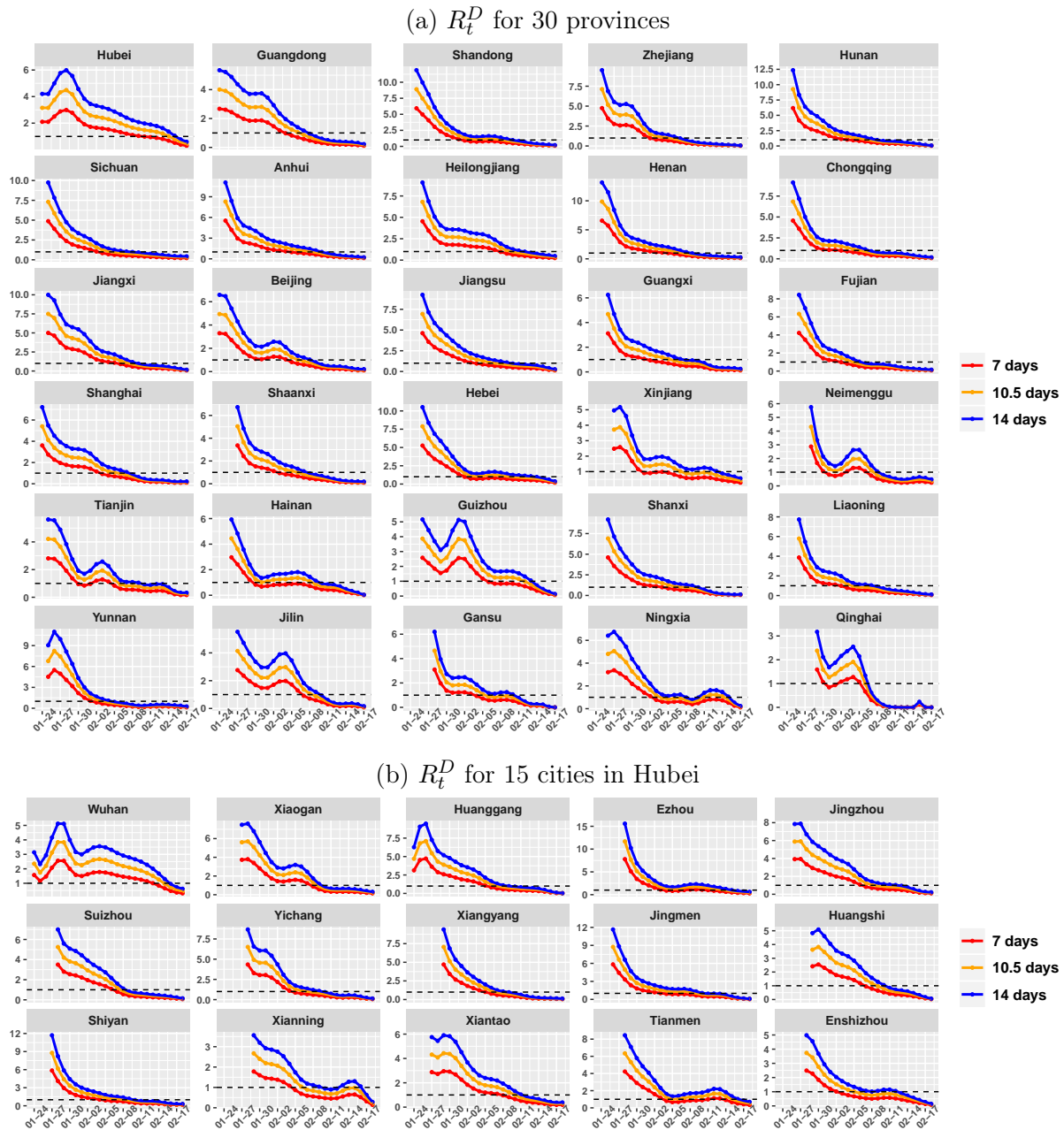


Figure 1: Time series of the reproduction number  $R_t^D$  at three infectious durations:  $D = 7$  (red), 10.5 (orange), 14 (blue), for the 30 mainland provinces (a) and the 15 cities in Hubei province (b) from Jan 21 to Feb 17 2020. The black horizontal line is the critical threshold level 1.

where  $\alpha(t)$  is the confirmation or diagnosis rate from  $E$  to  $I$ . The ODEs in (2) may specify the conditional means of the independent Poisson increments. However, like the SIR model, the states before  $I$  are not infectious.

The basic and the effective reproduction numbers (RN),  $R_0$  and  $R_t$ , are important notions in epidemiology as they quantify the reproduction ability of an epidemic at the start ( $R_0$ ) and

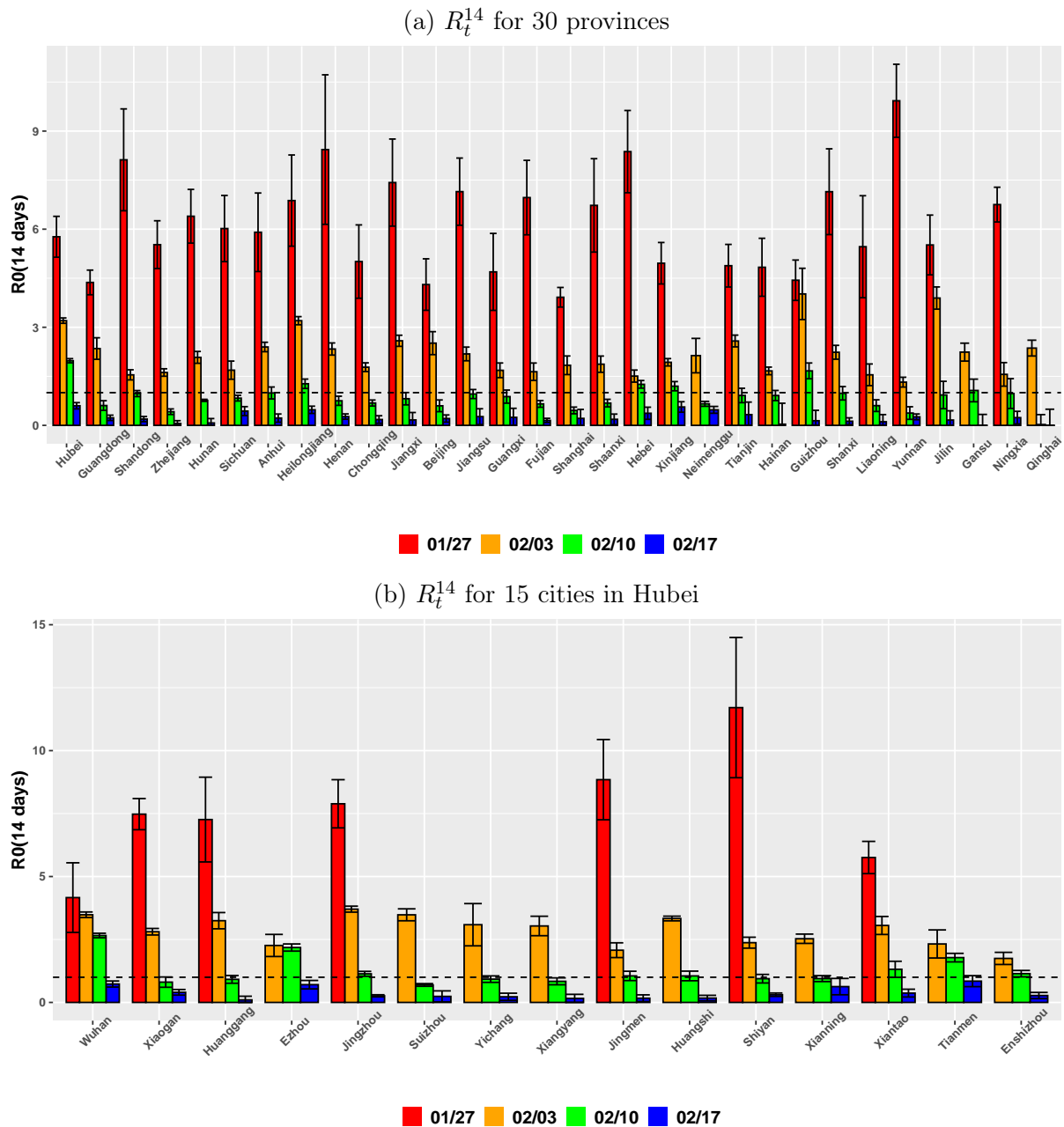


Figure 2: Elevated 95% confidence intervals (black) of the 14-day  $R_t$  for the 30 mainland provinces (a) and the 15 Hubei cities (b) on Jan 27 (red), Feb 3 (orange), Feb 10 2020 (green) and Feb 17 (blue). The black horizontal lines mark the critical threshold 1.

during ( $R_t$ ) an epidemic. For both the SIR and SEIR models,  $R_0 = R = \beta/\gamma$  (Hethcote, 2000). We will demonstrate that for the vSIR model,  $R_0 = \beta(0)/\gamma(0)$  and the effective RN  $R_t = \tilde{\beta}(t)/\gamma(t)$  at time  $t$  where  $\tilde{\beta}_t = \beta(t)s(t)$  and  $s(t) = S(t)/N$ . The susceptible rate  $s(t)$  is approximately 1 at the start of an epidemic. However,  $s(t) < 1$  has to be taken into account as the number of susceptibles declines.

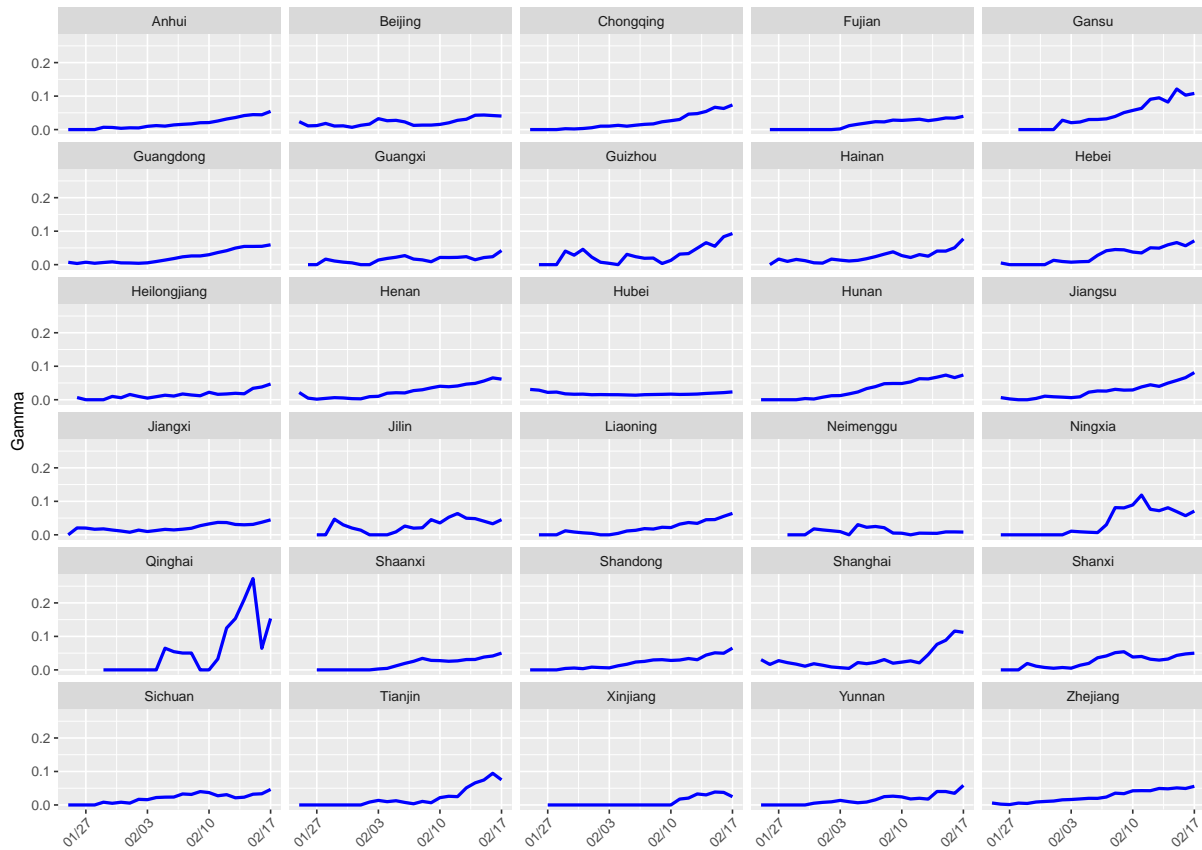


Figure 3: The estimated  $\hat{\gamma}(t)$  from the varying coefficient SIR model (1) for the data to Feb 17th 2020 for 30 provinces.

Let  $I(0)$  be the initial number of infected, according to the Poisson-vSIR model, at the start of the epidemic, conditioning on  $I(0)$ , in average  $I(1) = (1 - \gamma(0) + \beta(0))I(0) > I(0)$  if and only if  $1 - \gamma(0) + \beta(0) > 1$ , which is if and only if  $R_0 = \beta(0)/\gamma(0) > 1$ . In general, at time  $t$ , conditioning on  $I(t-1)$ , in average  $I(t) = (1 - \gamma(t-1) + \beta(t-1))I(t-1) > I(t-1)$  if and only if  $1 - \gamma(t-1) + \beta(t-1) > 1$ , which is the case if and only if the effective RN  $R_{t-1} > 1$ . Thus, indeed,  $R_t$  can track the trend of an epidemic being expanding or shrinking. A similar argument can be made under the vSEIR model.

## 4 Data

Daily records of infected, dead and recovered patients released by National Health Commission of China (NHCC) are obtained from the NHCC website, with the first confirmed record for Wuhan on December 8th, 2019, followed by 30 provinces in mainland China and 15 cities in Hubei province where Wuhan is the capital city. We did not consider data from Tibet due to very small number of cases. Due to severe under-reporting in the first 39 days of the epidemics in Wuhan and Hubei, we consider data from January 16th for Wuhan and Hubei. For other provinces and Hubei cities, the starting dates for data are those of first confirmed case, and the analysis date starts four days after to accommodate the estimation approach for the infectious rate  $\beta(t)$ . The



latest start data for analysis was January 29th for Qinghai province and three cities in Hubei province. The second last analysis starting date was January 28th with two provinces and five Hubei cities. Table A1 in the Supplementary Information (SI) provides the starting dates of the data records and analysis for each province and Hubei city.

To guide for the choice of the infectious duration  $D$  used when calculating the reproduction number, we consider two public data sources. The first one is obtained in *Shenzhen Government Online*, which contain datasets released by the Shenzhen Municipal Health Commission from January 19th to February 13th ([Shenzhen Municipal Affairs Service Data Administration, 2020](#)). One dataset provides information on the confirmed cases that include the time of onset, time of hospital admission, cause of illness and other information of 391 cases, consisting 188 males and 203 females. The admission time of these cases ranged from January 9th to February 11th. Another Shenzhen dataset reports the discharge times for 94 recovered cases, contained in the former dataset. The second data source comes from Shaoyang Municipal Health Committee ([Shaoyang Municipal Health Commission, 2020](#)) with a dataset of 100 confirmed cases released on February 14 that includes 48 male and 52 female patients with the onset dates ranging from January 12 to February 11.

## 5 Estimation and Confidence Intervals

The reported numbers of infected  $I(t)$  and removed cases  $R(t)$  are subject to measurement errors. To reduce the errors, we apply a three point moving average filter on the reported counts to obtain  $\bar{I}(t) = 0.3I(t-1) + 0.4I(t) + 0.3I(t+1)$  for  $2 \leq t \leq T-1$  where  $T$  is the latest time point of observation. In our analysis,  $T$  was February 20 of 2020. For  $t = 1$  or  $T$ , we apply two point averaging with  $7/10$  weight at  $t = 1$  or  $T$ , and  $3/10$  for  $t = 2$  or  $T-1$ . The above three-day moving window was chosen by assuming a possible one day reporting error, longer weight scheme may be applied, for instance a five day moving window may be used. The coefficient weights over the three days were made quite flat to reflect the nature of the reporting errors. Other weights can be considered. We expect the results would not be significantly affected by different weighting schemes. Apply the same filtering on the recovered process  $R(t)$  and obtain  $\bar{R}(t)$ . To simplify the notation, we denote the filtered data  $\bar{I}(t)$  and  $\bar{R}(t)$  as  $I(t)$  and  $R(t)$  respectively, wherever there is no confusion.

Hubei started to report the "clinically diagnosed" cases on February 12th (13th for city Xianning) which created spikes in the newly reported cases. We applied a one-side linear filter that re-distributes the spikes in the Hubei cities and Hubei to the previous 7 days with decreasing weights ranging from  $7/28$  to  $1/28$ .

Let  $N(t) = I(t) + R(t)$  denote the cumulative number of diagnosed cases, and let  $\Delta N(t) = N(t+1) - N(t)$  and  $\Delta R(t) = R(t+1) - R(t)$  denote the daily changes of  $N(t)$  and  $R(t)$ . Conditional on  $I(t)$ ,  $(\Delta N(t), \Delta R(t))$  are conditionally independent Poisson random variables. There is a slight confusion between  $N(t)$  and  $N$ , as the latter is used to denote the total population size.

We consider the likelihood for the Poisson-vSIR process framework for parameter estimation by treating  $\beta$  and  $\gamma$  as fixed and later we will relax it to allow they vary over a window of time  $t$ . Then,

$$\Delta N(t) \sim \text{Poisson} \{ \beta(t) S(t) I(t) / N \} \text{ and } \Delta R(t) \sim \text{Poisson} \{ \gamma(t) \}.$$

The likelihood function for  $(\Delta N(t), \Delta R(t))$  given  $I(t)$  and  $S(t)$  is

$$L(\beta, \gamma) = f_1(\Delta N(t) | I(t), S(t)) \times f_2(\Delta R(t) | I(t))$$

where  $f_1$  and  $f_2$  are the Poisson density functions with mean  $\beta(t)s(t)I(t)$  and  $\gamma(t)$ , respectively, where  $s(t) = S(t)/N$ . The log likelihood based on the increments at  $t$  is

$$l(\beta, \gamma) \propto -\beta(t)s(t)I(t) + \Delta N(t) \log\{\beta(t)s(t)I(t)\} - \gamma(t)I(t) + \Delta R(t) \log\{\gamma(t)I(t)\}.$$

As the population of each province/city is large and the number of total infected patients is still relatively small, the ratio  $S(t)/N$  appeared in (1) is very close to 1. By approximating  $s(t) = 1$ , the likelihood score equations are

$$\frac{\partial l}{\partial \beta} = -I(t) + \frac{\Delta N(t)}{\beta(t)} \quad (3)$$

$$\frac{\partial l}{\partial \gamma} = -I(t) + \frac{\Delta R(t)}{\gamma(t)} \quad (4)$$

It can be checked that the score functions have (approximate) zero means. The approximation of  $S(t)/N = 1$  is just to simplify the expression as everything carries through by using  $\tilde{\beta}(t) = \beta(t)s(t)$ .

While one can use the above likelihood based inference, an equivalent approach we use in our analysis is based on the (approximate) solution for  $I(t)$  via (5)

$$I(t) \approx I(t_1) \exp\{(\beta(t) - \gamma(t))(-t_1)\}, \quad (5)$$

for  $t_1 = t - w + 1, \dots, t$  and a window  $w > 0$  which satisfies  $w \rightarrow \infty$  and  $w/T \rightarrow 0$ . Here  $T$  is the total number of observational time for the processes. Take logarithm transform on (5),  $\log\{I(t)\} \approx \log\{I(t_1)\} + (\beta(t) - \gamma(t))(t - t_1)$ . We propose estimating  $\beta(t) - \gamma(t)$  by a local linear regression of  $\log\{I(t)\}$  on  $t - t_1$ . The above log-linear regression may be viewed as a version of the Poisson increment mean model by noting that  $\log\{I(t)\} - \log\{I(t_1)\} \approx \frac{I(t) - I(t_1)}{I(t_1)}$  which is approximately  $(\beta(t) - \gamma(t))(t - t_1)$  in the mean.

Let  $\widehat{\beta(t) - \gamma(t)}$  be the estimated slope from the local linear regression, and  $\widehat{\text{Var}}(\beta(t) - \gamma(t))$  be the estimated variance of  $\widehat{\beta(t) - \gamma(t)}$ . Their close form expressions are provided in Section S.1 in SI.

Let  $\Delta_\delta R_t = R_{t+\delta} - R_t$  for  $t = 1, \dots, T - \delta$ . From the second score equation (4), we estimate  $\gamma(t)$  by the local least square fitting of  $\Delta_\delta R_t$  on  $I(t)$  without intercept. Let  $\hat{\gamma}(t)$  and  $\widehat{\text{Var}}(\hat{\gamma}(t))$  be the estimator of  $\gamma$  and its corresponding estimated variance, respectively. Their expressions are provided in SI.

Then,  $\hat{\beta}(t) = \widehat{\beta(t) - \gamma(t)} + \hat{\gamma}(t)$  is the estimate for the varying coefficient  $\beta(t)$  in (1). The standard error of  $\hat{\beta}(t)$  can be obtained as

$$\text{SE}_\beta(t) = \left\{ \widehat{\text{Var}}(\widehat{\beta(t) - \gamma(t)}) + \widehat{\text{Var}}(\hat{\gamma}(t)) + 2\text{Cov}(\widehat{\beta(t) - \gamma(t)}, \hat{\gamma}(t)) \right\}^{1/2}.$$

The 95% confidence interval for  $\beta(t)$  can be constructed as

$$(\hat{\beta}(t) - 1.96\text{SE}_\beta(t), \hat{\beta}(t) + 1.96\text{SE}_\beta(t)).$$

In the implementation, we chose  $\delta = 1$  and  $w = 5$ . We had experimented other difference order  $\delta$  and found that  $\delta = 1$  offered better stability in the  $\gamma$  estimation. As the infection rates  $\beta(t)$  for all provinces declined quite rapidly over the study period, choosing  $w = 5$  reflect such a

change. One may use the cross-validation method to choose  $w$  and make it different for different provinces as well.

To assess the goodness of fitting, Figure S1 in SI shows the observed infected number  $I(t)$  versus the fitted values by the proposed varying coefficient SIR model for 30 provinces in China. It demonstrates the proposed method is well suitable for the dynamics of COVID-19 outbreak. Figure S2 in SI plots the estimated the effective reproductive number  $R_t^{14}$ , calculated as  $R_t^{14} = \hat{\beta}(t) \times 14$ , with its 95% confidence interval for 30 provinces in China.

## 6 Effective Reproduction Number

The effective reproduction number  $R_t$  is the most important parameter in determining the state of an epidemic. It measures the average number of infection made by an infectious person during the course of his/her being infectious. If  $R_t < 1 (> 1)$  for  $t$  larger than a  $t_0$ , then the epidemic will die down eventually (explode). There are two widely adopted definitions of  $R(t)$ . One is based on the average duration of infection of the disease, and the other is via the removal rate  $\gamma(t)$ .

At a date  $t$ , the effective reproduction number based on an average infectious duration  $D$  is  $R_t^D = \beta(t)D$  where  $\beta(t)$  is the daily infection rate at  $t$ . We do not adopt the version involving  $\gamma$ , the removal rate, since its estimation is highly volatile at the early stage of an epidemic. A general version of  $R(t)$  may be defined as  $\int_{t-D_1}^{t+D_2} \beta(u)du$  where positive  $D_1$  and  $D_2$  represent the infectious durations before and after diagnosis, respectively. The  $R_t^D$  given above can be viewed as an approximation by the Mean Value Theorem in calculus with  $D = D_1 + D_2$ .

Research works (Li et al., 2020; Guan et al., 2020; Chen et al., 2020) so far on COVID-19 have informed a range of duration for incubation, from onset of illness to diagnosis and then to hospitalization. The average incubation period from the three studies ranged from 3.0 to 5.2 days; the median duration from onset to diagnosis was 4 days (Guan et al., 2020); and the mean duration from onset to first medical visit and then to hospitalization were 4.6 and 9.1 days (Li et al., 2020), respectively. Based on a data sample of 391 cases from Shenzhen, the average incubation period was 4.46 (0.26) days and the average duration from onset to hospitalization were 3.9 (0.19) days, respectively, where standard error is reported in the parentheses. Another dataset of 100 confirmed cases in Shaoyang (Hunan Province) revealed the average durations from onset to diagnosis and from diagnosis to discharge were 5.67 (0.39) and 10.12 (0.43) days, respectively. There is a recent revelation (Guan et al., 2020) that asymptomatic patients can be infectious, which would certainly prolong the infectious duration.

There are much variation in the medical capability in timely diagnosis and hospitalization (thus quarantine) of the infected across the country. Thus, the infectious duration  $D$  would vary among the provinces and cities, and would change with respect to the stage of the epidemic as well.

Given the diverse range of infectious duration across the provinces and cities, in order to standardize and make the effective reproduction number  $R_t$  readily comparable, we calculated the  $R_t^D$  based on three levels of  $D$ : 7, 10.5 and 14 days, which represent three scenarios of responsiveness in diagnosing, hospitalization and hence quarantine of the infected. Calculation of the  $R_t$  at other duration can be made by inflating or deflating a  $R_t^D$  proportionally to reflect a local reality.

## 7 Reproductivity of COVID-19 in China

By calculating the time-varying infection rate function  $\beta(t)$ , we present in Figures 1 the time series of estimated  $R_t^D$  at the three levels of  $D$  for the 30+15 provinces/cities from late January to February 17th. Figure 2 displays four cross sectional  $R_t^{14}$  and their confidence intervals on January 27th, February 3rd, 10th and 17th, respectively.

Figure 1 reveals a monotone decreasing trend for almost all the provinces and cities with only exceptions for Hubei, Guizhou, Jilin, Neimenggu, and Qinghai. Even for those exceptional provinces, the recent trend is largely declining. The non-monotone pattern for non-Hubei provinces were largely due to relative small number of infected cases and waves of introduced infections. However, the one for Hubei and Wuhan suggests low data quality and in particularly under reporting and reporting delay. The epidemic statistics from Hubei and the city of Wuhan before January 21th were severely incomplete and with irregular patterns, as millions of people fled from Wuhan before the lockdown. This was the reason we start Hubei's analysis from January 21th.

The average  $R_t^{14}$  among the 27 provinces (with confirmed cases on and prior to January 23rd) was 6.14 (1.49), and 7.59 (2.38) for 7 of the 15 Hubei cities on January 27. These levels were comparable to the level of  $R$  (6.47) given in (Tang et al., 2020).

One week later on February 3rd,  $R_t^{14}$  was averaged at 2.18 (0.67) for the 30 provinces and 2.84 (0.59) for the 15 Hubei cities, indicating that cutting off Wuhan and other cities, and the start of wearing face masks and self isolation at home from January 23th had contributed to 64.5% and 62.6% reduction in the  $R_t^{14}$ . In the following week starting from February 4th, the average  $R_t^{14}$  came down to 0.86 (0.38) for the 30 provinces and 1.23 (0.55) for the 15 Hubei cities on February 10th, representing further 60.5% and 56.7% reductions, respectively, during the second week. This reflects the beneficial effects of the continued large scale self-isolation within the extended spring festival holiday period.

Table 1 provides the reproduction number  $R_t^D$  at the two durations on February 10th. It shows that 5 provinces and 5 Hubei cities'  $R_t^{14}$  were significantly above 1 (at 5% significance level). There are 14 provinces and 2 Hubei cities'  $R_t^{14}$  were significantly below 1, which were 1 and 1 more than those a day earlier on February 9th, and 9 and 2 more than those on February 8th, respectively. If we use the shorter  $D = 10.5$ , 22 provinces and 11 Hubei cities had been significantly below 1 for 1-7 consecutive days. These indicated that the reproduction number  $R_t$  has showed signs of crossing below the critical threshold 1 in increasing number of provinces and cities in Hubei around February 8-10. An updated Table 1 for February 17th are available in Table 2, which showed further improvement since February 10.

On February 17th,  $R_t^{14}$  of all provinces and cities under consideration have all been statistically significantly below 1, among which 22 provinces and 8 Hubei cities had been for at least seven consecutive days.

Given the significant decline in the reproduction numbers, it was time to discuss the turning point for COVID-19 for China. If a province or city's  $R_t^D$  started to be below 1 significantly (at 5% level), we would say the province or city have showed signs of the turning point. Given the uncertainty with the data records, especially those large variation in daily infected numbers coming out of Wuhan and Hubei, the turning point of the epidemic would be confirmed if  $R_t^D$  have been significantly below 1 for  $D_1$  days, where  $D_1$  is the period of infection before diagnosis, assuming all diagnosed cases can be quarantine immediately. Based on the results in (Li et al., 2020; Guan et al., 2020; Chen et al., 2020),  $D_1 = 7$  may be considered. Then, based on this criterion, some of the 30+15 provinces/cities had already reached the turning point on February

Table 2: The reproduction number  $R_t^D$  at two infectious durations: 10.5 and 14, for the 30 mainland provinces and 15 cities in Hubei province on February 17th. The symbols + (-) indicate that the  $R_t^{10.5}(R_t^{14})$  was significantly above (below) 1 at 5% level of statistical significance, and the numbers inside the square brackets were the consecutive days the  $R_t^{10.5}(R_t^{14})$  were significantly below 1.

Province/City	$R_t^{10.5}$	$R_t^{14}$	Province/City	$R_t^{10.5}$	$R_t^{14}$
Wuhan	0.55-[9]	0.73-[8]	Shanxi	0.09-[16]	0.13-[14]
Ezhou	0.52-[3]	0.69-[2]	Yichang	0.17-[17]	0.22-[14]
Tianmen	0.63-[9]	0.85	Yunnan	0.2-[21]	0.26-[19]
Hubei	0.43-[3]	0.56-[2]	Anhui	0.16-[8]	0.21-[7]
Sichuan	0.33-[19]	0.44-[16]	Xianning	0.47-[8]	0.63-[8]
Xiantao	0.28-[14]	0.37-[13]	Suizhou	0.18-[16]	0.25-[16]
Tianjin	0.25-[11]	0.33-[10]	Shandong	0.14-[12]	0.19-[9]
Heilongjiang	0.34-[7]	0.46-[5]	Guangdong	0.17-[10]	0.23-[9]
Shiyan	0.23-[15]	0.31-[14]	Xiangyang	0.12-[17]	0.16-[15]
Neimenggu	0.36-[17]	0.47-[16]	Zhejiang	0.04-[18]	0.06-[17]
Xiaogan	0.28-[14]	0.37-[13]	Enshizhou	0.2-[12]	0.27-[4]
Xinjiang	0.42-[12]	0.56-[10]	Huanggang	0.07-[10]	0.09-[7]
Beijing	0.13-[11]	0.18-[9]	Guizhou	0.11-[4]	0.15-[3]
Jingzhou	0.19-[8]	0.25-[3]	Jingmen	0.12-[8]	0.16-[4]
Shaanxi	0.14-[17]	0.18-[16]	Gansu	0-[7]	0-[6]
Chongqing	0.13-[11]	0.17-[10]	Hainan	0.03-[8]	0.04-[7]
Henan	0.2-[9]	0.26-[8]	Hunan	0.02-[10]	0.07-[9]
Hebei	0.26-[5]	0.35-[3]	Ningxia	0.18-[9]	0.24-[9]
Guangxi	0.17-[12]	0.23-[7]	Huangshi	0.13-[8]	0.17-[7]
Shanghai	0.16-[18]	0.21-[16]	Jiangxi	0.12-[9]	0.16-[7]
Jiangsu	0.2-[10]	0.26-[6]	Liaoning	0.08-[20]	0.11-[16]
Jilin	0.11-[7]	0.15-[7]	Qinghai	0-[4]	0-[4]
Fujian	0.11-[13]	0.15-[11]			

17, and more would follow in the coming days according to latest Table 2

## 8 Prediction for Infection Rate and State Variables

As  $R_t^D = \beta(t)D$ , predicting  $\beta(t)$  is equivalent to predicting  $R_t^D$ . From Figure 1 and Figure S2 in SI, we see that the overall trends of  $\beta(t)$  is decreasing. But the rate of decreasing became smaller as time travels. As the nonparametric estimates for  $\beta(t)$  are harder to be extended beyond the data range, we consider fitting the estimated  $\beta(t)$  with a parametric model, and use the latter for projection to the future. Specifically, to reflect the declining trend of  $\beta(t)$ , we consider the reciprocal regression

$$\beta(t) = \frac{b}{t^\eta - a} + e_t \quad (6)$$

with error  $e_t$  and unknown parameters  $a$ ,  $b$  and  $\eta$ . The parameters  $a$  and  $b$  are estimated by minimizing the sum-of-squares distance between the last 14 day estimates  $\hat{\beta}(t)$  and the fitted

values for a given  $\eta$ , and then the optimal  $\eta$  is chosen to be the one that gives the minimum mean square error over a set of candidate values from 0.5 to 5 with 0.1 increment. Let  $\tilde{a}$ ,  $\tilde{b}$  and  $\tilde{\eta}$  be the estimated parameters, and  $\check{\beta}(t) = \tilde{b}/(t^{\tilde{\eta}} - \tilde{a})$  be the fitted function. Figure S3 in SI shows the reciprocal model fits  $\hat{\beta}(t)$  quite well for most of the provinces, especially those with large number of infected cases.

With the fitted  $\check{\beta}(t)$ , we project  $\{S(t), I(t), R(t)\}$  via the ODEs

$$\begin{aligned}\frac{d\hat{S}(t)}{dt} &= -\check{\beta}(t)\hat{I}(t)\frac{\hat{S}(t)}{N}, \\ \frac{d\hat{I}(t)}{dt} &= \check{\beta}(t)\hat{I}(t)\frac{\hat{S}(t)}{N} - \hat{\gamma}_T\hat{I}(t), \\ \frac{d\hat{R}(t)}{dt} &= \hat{\gamma}_T\hat{I}(t).\end{aligned}\tag{7}$$

where  $\hat{\gamma}_T$  is the estimated recovery rate at time  $t$  using the last five days' data. With the observed  $\{S(T), I(T), R(T)\}$  at the current time  $t$  as the initial values, numerical solutions  $\{(\hat{S}(t), \hat{I}(t), \hat{R}(t)) : T \leq t < \infty\}$  for the system (7) could be obtained using the Euler method. Then, the end time of the epidemic can be predicted as  $t_{\text{end}} = \min\{t : \hat{I}(t) < 1\}$ , and the estimated final infected number is  $\hat{N}_{\text{final}} = \hat{R}(t_{\text{end}}) + \hat{I}(t_{\text{end}})$ .

To conduct statistical inference for the epidemic predictions, we use the bootstrap method. In particular, we generate parametric bootstrap resampled processes based on the vSIR model which facilitate the construction of prediction intervals. We regard that the increments of  $S(t)$  and  $R(t)$  follow the Poisson processes (Bretó et al., 2009) over time as

$$-\Delta S(t) \sim \text{Poisson}\{\beta(t)S(t)I(t)/N\} \text{ and } \Delta R(t) \sim \text{Poisson}\{\gamma I(t)\}.$$

With the estimated  $\hat{\gamma}$  and  $\hat{\beta}(t)$ , we generate bootstrap samples  $\{(S^{(b)}(t), I^{(b)}(t), R^{(b)}(t))\}_{t=1}^T$  of the original process for  $b = 1, 2, \dots, B$ .

For each bootstrap resampled  $\{(S^{(b)}(t), I^{(b)}(t), R^{(b)}(t))\}_{t=1}^T$ , we obtain the estimates  $\beta_{\star}^b(t)$  and  $\gamma_{\star}^b$  for  $\beta(t)$  and  $\gamma$  in the same way as for the original sample. Let  $\bar{\beta}_{\star}(t) = \sum_{b=1}^B \beta_{\star}^b(t)/B$  and  $\bar{\gamma}_{\star} = \sum_{b=1}^B \gamma_{\star}^b/B$  be the average of the bootstrap estimates. We employ the bias corrected bootstrap estimates for  $\beta(t)$  and  $\gamma$  as

$$\hat{\beta}^b(t) = \hat{\beta}(t) + (\beta_{\star}^b(t) - \bar{\beta}_{\star}(t)) \quad \text{and} \quad \hat{\gamma}^b = \hat{\gamma} + (\gamma_{\star}^b - \bar{\gamma}_{\star})$$

for  $b = 1, 2, \dots, B$ . We then use the reciprocal model (6) to project the future path of  $\hat{\beta}^b(t)$ , and use the numerical solution of the vSIR ODEs to predict the end time and the accumulative number of final infected cases as we described in section 3.4. Let the bootstrap estimates for the peak time be  $\{t_{\text{end}}^b\}_{b=1}^B$ . The 95% prediction interval for the peak time is constructed as the 2.5% and 97.5% quantiles of  $\{t_{\text{end}}^b\}_{b=1}^B$ . Similar bootstrap prediction intervals can be constructed for the final accumulative infection number  $N_{\text{final}}$  of the epidemic.

## 9 Prediction Results

Based on the estimated  $\beta(t)$  over time, we predict COVID-19's future trajectories as solutions to the vSIR model. We used data up to February 19 2020 for the prediction under three scenarios for the recovery rate  $\gamma$ . One uses the empirical estimate based on data to February 19th. As an

effective cure for the virus has not been found, the estimated recovery rates are quite low. Among the provinces with more than 100 infections on February 19, Jiangsu had the highest recovery rate 0.08, followed by Jiangxi, Hebei, Shanghai, Shanxi, Chongqing, Henan (0.07). Hubei, the province at the center of the epidemic, is 0.025. The other scenarios was to choose  $\gamma = 1/14$  and  $\gamma = 0.1$ , which mean the average removal time from diagnosis was 14 and 10 days, respectively, representing improvement in the treatment for COVID-19 patients as time progressed.

Tables 3 presents the 95% prediction intervals for the end times of the epidemic and the cumulative number of infected at the ending. The trajectories of  $I(t)$  of the proposed vSIR model are presented in Figure 4 under the three scenarios of the recovery rate. The predicted infection

Table 3: The 95% prediction intervals for the ending times, and the final accumulative number of infected cases of COVID-19 epidemic in the 30 provinces based on data to Feb 19 2020 with  $\gamma = 0.1$ . The last column lists the total infected cases ( $I(t) + R(t)$ ) as Feb 19, 2020.

Province	Peak time	Ending time	$\hat{N}_{\text{final}}$	Current
Hubei	2/20 - 2/22	6/20 - 6/21	73857 - 74596	62322
Guangdong	2/9 - 2/9	4/27 - 4/29	1368 - 1412	1347
Zhejiang	2/7 - 2/7	4/26 - 4/27	1225 - 1245	1195
Beijing	2/11 - 2/20	4/17 - 4/20	416 - 436	397
Chongqing	2/10 - 2/10	4/18 - 4/21	581 - 600	565
Hunan	2/10 - 2/10	4/21 - 4/23	1028 - 1046	1021
Guangxi	2/12 - 2/12	4/11 - 4/15	254 - 271	248
Shanghai	2/9 - 2/9	4/12 - 4/16	345 - 365	336
Jiangxi	2/14 - 2/14	4/23 - 4/25	969 - 994	955
Sichuan	2/14 - 2/23	4/25 - 4/28	589 - 619	525
Shandong	2/23 - 3/20	4/19 - 4/21	567 - 584	553
Anhui	2/11 - 2/24	4/26 - 4/28	1044 - 1068	1006
Fujian	2/10 - 2/10	4/13 - 4/16	306 - 320	299
Henan	2/9 - 2/9	4/29 - 5/1	1358 - 1387	1283
Jiangsu	2/15 - 2/15	4/20 - 4/23	662 - 687	640
Hainan	2/19 - 3/9	4/6 - 4/9	174 - 183	168
Tianjin	2/19 - 3/7	4/7 - 4/14	141 - 159	132
Yunnan	2/13 - 2/13	4/7 - 4/11	174 - 187	174
Shaanxi	2/10 - 2/10	4/12 - 4/16	262 - 276	250
Heilongjiang	2/19 - 2/26	4/23 - 4/26	519 - 554	479
Liaoning	2/9 - 2/9	3/31 - 4/3	122 - 126	122
Guizhou	2/12 - 2/12	4/4 - 4/10	150 - 165	147
Jilin	2/9 - 2/9	3/30 - 4/4	93 - 101	92
Ningxia	2/13 - 2/13	3/18 - 3/26	65 - 73	71
Hebei	2/19 - 3/18	4/12 - 4/16	319 - 336	312
Gansu	2/9 - 2/9	3/20 - 3/26	92 - 96	92
Xinjiang	2/19 - 2/29	3/31 - 4/9	78 - 96	78
Shanxi	2/10 - 2/10	4/1 - 4/5	134 - 142	134
Neimenggu	2/14 - 2/14	4/2 - 4/11	78 - 98	76
Qinghai	2/4 - 2/4	2/23 - 3/6	19 - 20	19

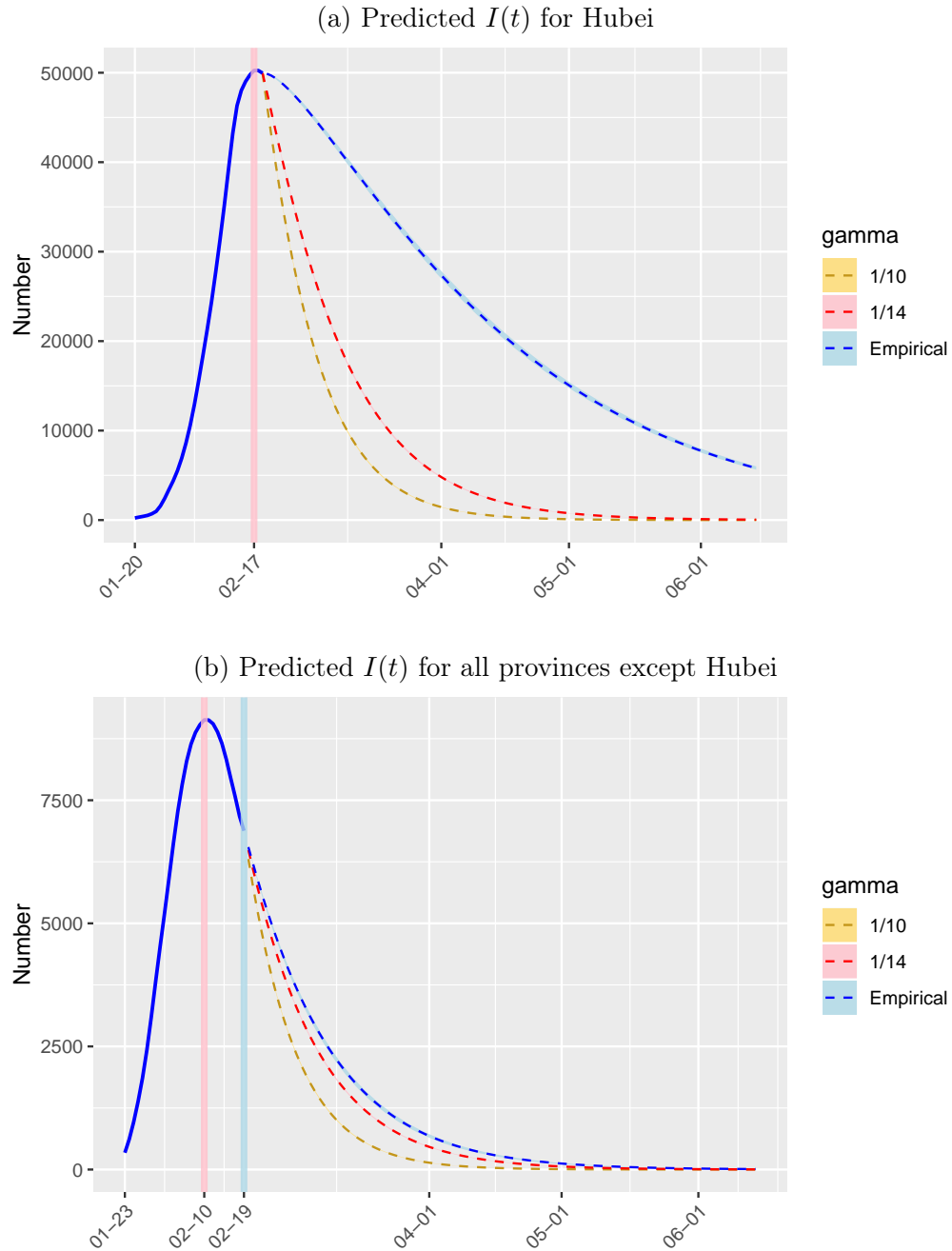


Figure 4: Predicted number of infected cases  $I(t)$  with 95% prediction interval for Hubei Province in panel (a) and all other provinces combined except Hubei in panel (b). The grey vertical line indicates the current date of observation; the blue solid line plots the observed  $I(t)$  before Feb 19th; the blue dashed line gives the predicted  $I(t)$  with 95% prediction interval (blue shaded area) with the estimated  $\hat{\gamma}_T$ ; the pink vertical line indicates the peak date of  $I(t)$ ; the orange and red dashed line gives the predicted  $I(t)$  with 95% prediction interval (shaded area) with fixed recovery rate  $\gamma = 0.1$  and  $\gamma = 1/14$  respectively.



number  $\hat{I}(t)$  is within 10% deviation from its observed value based on data up to Feb 19th, see Table A2 in SI for the detailed prediction error.

From the trajectory of the vSIR model in Figure 4, for the non-Hubei provinces, the number of infected would quickly decrease in late February and March with very few cases left in April under all the three scenarios. Some provinces with few number of total infected cases may end as early as March (Qinghai, Jilin, Gansu, Ningxia). For Hubei, with a higher recovery rate of 0.1, the duration of the epidemic would be shorten substantially. The ending time for Hubei is around June 20 2020 with total number of infection in the range 73,857–74,596. This shows that improving the recovery rate is an efficient way to end the COVID-19 infection early given the current decreasing trend of  $\beta(t)$ , as it leads to the reduction of the infectious duration.

## 10 Discussion

The implications of China's experience in combating COVID-19 to other countries facing the epidemic are two folds. One is to reduce the person-to-person contact rate by self isolation and curtailing population movement; another is to reduce the transmission probability by wearing protective wears when a contact has to be made.

The eventual control of COVID-19 is rested on if the existing control measures can be continued further for a period of time. The biggest challenges that can jeopardize the great effort from late January are from the impatient populations eager to get out of the self-isolation driven by either economic needs (migrant workers eager to coming back to cities for income) or people trying to escape from the self isolation encouraged by the declining infections in the last two weeks. In any case, the vSIR model and its statistical estimation and inference can be used to model and the assess the COVID-19 epidemics in other countries.

## Supplementary Materials

The data and R code are publicly available with explanations on various segments of the procedures that we have proposed in this paper; see <https://github.com/sun-haoxuan/vSIR>. The supplementary information referred to in the paper can be found on the *Journal of Data Science* website.

## References

- Ball F, Clancy D (1993). The final size and severity of a generalised stochastic multitype epidemic model. *Advances in Applied Probability*, 25(4): 721–736.
- Becker NG (1977). On a general stochastic epidemic model. *Theoretical Population Biology*, 11(1): 23–36.
- Becker NG, Britton T (1999). Statistical studies of infectious disease incidence. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 61(2): 287–307.
- Bretó C, He D, Ionides EL, King AA, et al. (2009). Time series analysis via mechanistic models. *The Annals of Applied Statistics*, 3(1): 319–348.
- Chen N, Zhou M, Dong X, Qu J, Gong F, Han Y, et al. (2020). Epidemiological and clinical characteristics of 99 cases of 2019 novel coronavirus pneumonia in Wuhan, China: A descriptive study. *The Lancet*, 395: 507–513.

- Cleveland WS, Devlin SJ (1988). Locally weighted regression: An approach to regression analysis by local fitting. *Journal of the American Statistical Association*, 83(403): 596–610.
- Guan WJ, Ni ZY, Hu Y, Liang WH, Ou CQ, He JX, et al. (2020). Clinical characteristics of coronavirus disease 2019 in China. *New England Journal of Medicine*, 382: 1708–1720.
- Hethcote HW (2000). The mathematics of infectious diseases. *SIAM review*, 42(4): 599–653.
- Kermack WO, McKendrick AG (1927). A contribution to the mathematical theory of epidemics. *Proceedings of the Royal Society of London, Series A*, 115(772): 700–721.
- Li Q, Guan X, Wu P, Wang X, Zhou L, Tong Y, et al. (2020). Early transmission dynamics in Wuhan, China, of novel coronavirus-infected pneumonia. *New England Journal of Medicine*, 382(13): 1199–1207.
- Read JM, Bridgen JR, Cummings DA, Ho A, Jewell CP (2020). Novel coronavirus 2019-nCoV: Early estimation of epidemiological parameters and epidemic predictions. MedRxiv preprint: <https://doi.org/10.1101/2020.01.23.20018549>.
- Shaoyang Municipal Health Commission (2020). Dynamic information on prevention and control of new coronavirus-infected pneumonia in Shaoyang. Available at <https://wjw.shaoyang.gov.cn/wjw/zyxw/202002/c49df53092784c85aaac769149f30265.shtml>. [Accessed February 14, 2020].
- Shenzhen Municipal Affairs Service Data Administration (2020). Latest data. <https://opendata.sz.gov.cn/data/dataSet/toDataSet>. [Accessed February 14, 2020].
- Tang B, Wang X, Li Q, Bragazzi NL, Tang S, Xiao Y, et al. (2020). Estimation of the transmission risk of the 2019-nCoV and its implication for public health interventions. *Journal of Clinical Medicine*, 9(2): 462.
- World Health Organization (2020). Situation Report — 26. Available at <https://www.who.int/docs/default-source/coronaviruse/situation-reports/20200215-sitrep-26-covid-19.pdf>. [Accessed February 16, 2020].
- Wu JT, Leung K, Leung GM (2020). Nowcasting and forecasting the potential domestic and international spread of the 2019-nCoV outbreak originating in Wuhan, China: A modelling study. *The Lancet*, 395(10225): 689–697.
- Yip PS, Chen Q (1998). Statistical inference for a multitype epidemic model. *Journal of Statistical Planning and Inference*, 71(1-2): 229–244.

## Discussion of “Tracking reproductivity of COVID-19 epidemic in China with varying coefficient SIR model”

YUKANG JIANG<sup>1</sup>, JIANBIN TAN<sup>1</sup>, TING TIAN<sup>1</sup>, AND XUEQIN WANG <sup>\*2</sup>

<sup>1</sup>*School of Mathematics, Sun Yat-sen University, Guangzhou, Guangdong, China*

<sup>2</sup>*School of Management, University of Science and Technology of China, Hefei, Anhui, China*

Sun et al. proposed a vSIR model with both time-varying transmission rate and removal rate. Based on the traditional SIR (Kermack and McKendrick, 1927) and SEIR (Hethcote, 2000) models, the authors introduced some random mechanisms for the daily changes of infected ( $I$ ) and removal ( $R$ ) individuals by assuming the Poisson increments:  $\Delta I(t)$  and  $\Delta R(t)$ . Also, the effective reproduction numbers calculated by the vSIR model were time-varying with corresponding 95% confidence intervals, which could effectively reflect the intervention policies of COVID-19 in different regions. In our discussion, we focus on the following points.

**The spreads of COVID-19** The spreads of COVID-19 in different regions of China were accurately estimated. From January 27 to February 17, 2020, the effective reproduction numbers of COVID-19 in 30 provinces and 15 cities in Hubei Province decreased significantly. The effective reproduction numbers had dropped below 1 from February 10 to 17, 2020 in most provinces, indicating that the epidemic situations in most areas in China had been controlled after the implementation of a series of interventions and control measures, which was similar to the conclusion of Tan et al. (2020).

The authors also predicted the ending times of the COVID-19 epidemic through the model by considering the numbers of cumulative infected cases. By setting the recovery rate to 0.1, their model estimated that the number of infected individuals would be significantly reduced during March in China. They obtained the conclusion that the epidemic of COVID-19 outside Hubei province would be approaching to an end in April. These estimations gave rise to a realistic predicted trend of the epidemic, which were close to those of Cui and Hu (2020).

**Estimation of  $R_t^{14}$**  The  $R_t^{14}$  was overestimated. The basic reproduction number ( $R_0$ ) can be considered as the expected number of cases directly infected by one primary case in a population where all individuals are susceptible to infection (Anderson and May, 1992). According to  $R_0$ s of COVID-19 reported by Liu et al. (2020),  $R_0$ s were between 1.4 and 6.49, with an average of 3.28 and a median of 2.79. For the 14-day time-varying effective reproduction number  $R_t^{14}$  defined in the vSIR model,  $R_t^{14}$  of 7 cities in Hubei province reached 7.59 on January 27. Since January 23, 2020, the Chinese government has taken stringent measures such as traffic blockade to curb the epidemic. The  $R_0$  of COVID-19 could be similar to the time-varying effective reproduction number in Hubei province in the early period, or even higher than  $R_t^{14}$  of Hubei province on January 27 as reported by Tan et al. (2020). Thus, the value of  $R_t^{14}$  was relatively high in the early stage of COVID-19 in Hubei province, which may be overestimated.

**Factors Affecting the Infection Rate** This article did not take into account the infectivity of the infected individuals or the population of migrants between different areas and some parameters are needed to be given in advance. There could be more considerations in the vSIR model.

---

\*Corresponding author. Email: wangxq20@ustc.edu.cn.

For example, the confirmed cases could not be infectious once they have been quarantined, while the infected individuals who should be considered as infectious hosts during the incubation period (Tan et al., 2020; Yang et al., 2020), which is not reflected in the SIR or SEIR model. At the same time, if the method is extended to various countries, the impact of population mobility in different regions could be considered (Yang et al., 2020; Gilbert et al., 2020). Besides, for the predictions from the model in the paper, the parameters in the vSIR model were needed to be specified beforehand. When the model is used to different countries and regions, it may be necessary to adjust the corresponding parameters subjectively to get better prediction results.

**Summaries** The authors proposed an adjusted SIR model: vSIR model, which can calculate effective reproduction numbers of different regions over time, and make predictions for the future trend. It is effective to evaluate the early epidemic situations in different regions. However, the paper overestimated the value of  $R_t^{14}$  in the early stage of the COVID-19 in China and the parameters in the vSIR model were needed to be specified in advance. Furthermore, if the infectivity of infected individuals during the incubation period and the mobility of the population are incorporated into the model, it would make the model more realistic.

## Acknowledgment

Wang's work was supported in part by the National Natural Science Foundation of China (Grant No.11771462).

## References

- Anderson RM, May RM (1992). *Infectious Diseases of Humans: Dynamics and Control*. Oxford university press.
- Cui H, Hu T (2020). Nonlinear regression in COVID-19 forecasting. *Scientia Sinica Mathematica*. Forthcoming, <https://doi.org/10.1360/SSM-2020-0055>.
- Gilbert M, Pullano G, Pinotti F, Valdano E, Poletto C, Boëlle PY, et al. (2020). Preparedness and vulnerability of African countries against importations of COVID-19: A modelling study. *The Lancet*, 395(10227): 871–877.
- Hethcote HW (2000). The mathematics of infectious diseases. *SIAM Review*, 42(4): 599–653.
- Kermack WO, McKendrick AG (1927). A contribution to the mathematical theory of epidemics. *Proceedings of the Royal Society of London, Series A*, 115(772): 700–721.
- Liu Y, Gayle AA, Wilder-Smith A, Rocklöv J (2020). The reproductive number of COVID-19 is higher compared to SARS coronavirus. *Journal of Travel Medicine*, 27(2): taa021.
- Tan J, Jiang Y, Tian T, Wang X (2020). P-SIHR probabilistic graphical model: An applicable model of COVID-19 in estimating the number of infectious individuals. *Acta Mathematicae Applicatae Sinica*, 43(2): 365–382.
- Yang Z, Zeng Z, Wang K, Wong SS, Liang W, Zanin M, et al. (2020). Modified SEIR and AI prediction of the epidemics trend of COVID-19 in China under public health interventions. *Journal of Thoracic Disease*, 12(3): 165–174.

## Discussion of “Tracking reproductivity of COVID-19 epidemic in China with varying coefficient SIR model”

LU TANG\*<sup>1</sup>

<sup>1</sup>*Department of Biostatistics, Graduate School of Public Health, University of Pittsburgh, Pittsburgh, Pennsylvania, USA*

The paper by Sun et al. is a thoughtful and important work for tracing the transmission of COVID-19 over the course of outbreak in China. As COVID-19 has been spreading globally as a pandemic, the lessons obtained from disease mitigation measures in China become more valuable to other countries and to future infectious diseases.

The authors adopted an epidemiological model, SIR model (Kermack and McKendrick, 1927), of infectious diseases, and made a few key modifications. First, a stochastic version of the SIR model, Poisson-vSIR model, is proposed to model the incremental numbers of infected and recovered cases by Poisson distributions, which facilitates statistical inference; more importantly, the infectious rate is allowed to vary over time and estimated by the reciprocal regression. The proposed model is applied to cases in provinces and cities of China to calculate the time-varying infection rate function. Results show that infectious rates significantly decreased after reducing the person-to-person contact and curtailing population movement, a valuable lesson for other countries facing surging COVID-19 cases and deaths. In what follows, I will discuss about the model, the preprocessing of observed data, the effective reproduction number, and the predictions of model.

The proposed Poisson-vSIR model is based on the Euler method, a numeric method for ordinary differential equations (Lapidus and Seinfeld, 1971, ODEs), discretizing cases arising from the continuous time scale to cases modelled at the discretized time scale in days. Recognizing this numeric step would help appreciate the connection between the ODEs-based SIR model and the difference equations (DEs)-based Poisson-vSIR model. In the same time, I am curious if this numeric approximation would cause some discrepancy between two models in terms of the dynamic of diseases and if higher order of numeric approximation of ODEs, like the Runge–Kutta method, would make any difference.

A key component of Poisson-vSIR model is the time-varying infection rate function, pre-specified as a parametric form possibly due to identifiability concerns. There likely exist many other parametric forms which may not be distinguishable from the current one, given the limited dataset. What interests me is the observation that the continuous infection rate function fits well for most of the provinces, even when decisions of reducing the person-to-person contact, e.g. lockdown of city or “shelter in place”, were commonly effective shortly after the announcement, and behaviors of person-to-person contact changed abruptly. Do the results of this analysis suggest the opposite, that behaviors of person-to-person contact change gradually after the abrupt lockdown? Alternatively, could it be possible that the sudden changes of infection rate is smoothed out by data preprocessing. The numbers of infected and recovered cases are subjected to underreporting bias, availability of test kits, and changes of diagnosis and report criteria, among other factors. For example, 15,152 new cases were reported of February 12 2020 in China, about 600% surge over the preceding day, largely due to the changes of criteria how cases are diagnosed and reported. Without adjusting such “outliers”, the results from the Poisson-vSIR model would be misleading. To deal with it, moving average filter was applied. While moving average will

---

\*Email: lutang@pitt.edu.

filter outliers, it might also smooth out the short-term exponential growth trend or some change points due to intervention policies. Additional sensitivity analysis would be of interest.

The main story of paper is built around the effective reproduction number. The authors considered two possible forms:  $R_t = \tilde{\beta}(t)/\gamma(t)$ , where  $\tilde{\beta}(t)$  is the filtered infectious rate and  $\gamma(t)$  is the removal rate; and  $R_t^D = \beta(t)D$ , where  $\beta(t)$  is the unfiltered infectious rate and  $D$  is an average infectious duration, which can be viewed as a special case with  $\gamma(t)$  remains constant over time. Authors preferred  $R_t^D$  and pointed out that  $\gamma(t)$  involved in  $R_t$  is highly volatile at the early stage of COVID-19, which makes it critical to choose the proper plug-in value of average infectious duration  $D$ . Even so, the variability of such  $D$  could also make a difference in constructing confidence intervals for  $R_t^D$ . As COVID-19 comes to an end in China, it may be worthwhile revisiting the point estimation and variance of  $D$ , and updating the estimation of  $R_t^D$  accordingly. In addition, how volatile the  $\beta(t)$  would be at the early stage of COVID-19 and if changes of  $R_t^D$  over time are statistically significant?

Policies of public mitigation rely on the predicted trajectory of infectious disease more than ever. For example, in order to justify the prolonged “stay at home” policy, members of the US Coronavirus Task Force explained the future outlook of COVID-19 in US based on the prediction of epidemiological models during a press briefing on March 30, 2020. To examine the promise of data-driven policy making, it would be of interest to evaluate if the observed trajectories of COVID-19 in provinces or cities of China follow the predicted ones by the proposed model closely and what lessons have been gained from predicting trajectories of infectious disease.

After fighting through the COVID-19 pandemic with blood, toil, and tears, it is the time to reflect. The tools developed by Sun et al. and lessons offered would be very valuable for a swift response to the next pandemic. More such efforts are desired.

## Acknowledgments

This work was supported by research startup fund of the Graduate School of Public Health, University of Pittsburgh, PA, USA.

## References

- Kermack WO, McKendrick AG (1927). A contribution to the mathematical theory of epidemics. *Proceedings of the Royal Society of London, Series A*, 115(772): 700–721.
- Lapidus L, Seinfeld JH (1971). *Numerical Solution of Ordinary Differential Equations*. Academic Press, New York.

## Discussion of “Tracking reproductivity of COVID-19 epidemic in China with varying coefficient SIR model”

LILI WANG<sup>1</sup>, FEI WANG<sup>2</sup>, YIWANG ZHOU<sup>1</sup>, AND PETER X.K. SONG<sup>\*1</sup>

<sup>1</sup>*Department of Biostatistics, University of Michigan, Ann Arbor, Michigan, USA*

<sup>2</sup>*Data Science Team, CarGurus, Cambridge, Massachusetts, USA*

We congratulate the authors on an interesting contribution to the Susceptible-Infection-Removed (SIR) modeling. In this paper the authors have focused on the development of a fast and flexible framework to incorporate a time-varying effective reproduction number ( $\mathcal{R}_t$ ) into the SIR model. Through estimation and extrapolation of such a time-varying effective reproduction number ( $\mathcal{R}_t$ ), this new model can predict time-varying reproduction profiles critical to the understanding of covid-19 disease evolution in China. The use of a reciprocal regression technique enables forecast the number of future infectious cases with a certain confidence level.

Let the numbers of individuals in the compartments of susceptible, infected and removed be, respectively,  $S_t$ ,  $I_t$  and  $R_t$ , which are changing over time  $t$ . Let  $N_t = I_t + R_t$  be the cumulative number of infected cases and  $N = S_t + I_t + R_t$  be the target population under investigation, which is fixed constant over time. To reduce noise, the authors applied a local smoothing by the moving average of three data points on the three compartment time series. The authors proposed a local linear fitting approach to estimating the nonlinear transmission rate  $\beta_t$  and removal rate  $\gamma_t$  at each target time  $t$  conditional on  $I_t$  via the following two local linear models at a target time  $t_1$ :  $\log(I_t) \sim \log(I_{t_1}) + (\beta_t - \gamma_t)(t - t_1)$ , and  $R_{t_1+\delta} - R_{t_1} \sim \gamma_t I_{t_1}$ , where  $\delta > 0$  is set to be 1. Because the number of removed cases, i.e. a total number of deaths and recovered cases, are often collected and reported mostly by hospitals several weeks after virus contract and infection, the observed series  $R_t$  is typically delayed and incomplete, and thus there is a concern with the reliability of such data in a small local time window used in the estimation. The data quality issue may result in unstable local estimation of the transmission rate  $\gamma_t$  at a time point. The authors suggested using  $\mathcal{R}_t^D = \beta_t D$  to estimate the time-varying effective reproduction number, where  $D$  is a prefixed infectious duration varying from 7 to 14 days.

For the prediction purpose, the authors considered a parsimonious nonlinear model that enables extrapolation of the estimated  $\beta_t$  into the future time. The transmission rate is assumed to satisfy the reciprocal regression model

$$\beta_t = \frac{b}{t^\eta - a} + e_t,$$

where  $a$ ,  $b$ , and  $\eta$  are the parameters to be estimated and  $e_t$  is the error. In addition, both projected  $N_t$  and  $R_t$  are obtained following the SIR ordinary differential equations using Euler method. The bias correction and estimating inference were achieved via a bootstrap method. In particular, we notice that the confidence intervals provided in their Figure 5 were very narrow, which does not seem to capture much certainty from the infection dynamic system well. This underestimation of uncertainty might result from maximizing a conditional likelihood (given  $I_t$ ) in their equation (3) and data pre-processing via an average filtering procedure.

To gain some insights of the above reciprocal regression model for the transmission rate, we implemented this time-varying  $\beta_t$  in an analysis of the COVID19 data of Hubei, China using

---

\*Corresponding author. Email: pxsong@umich.edu.

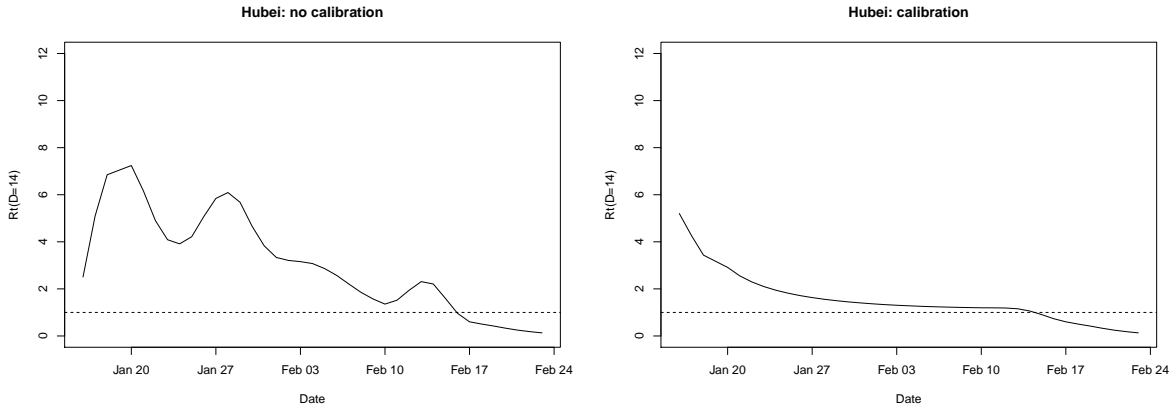


Figure 1: The time-varying effective reproduction number  $\mathcal{R}_t^{14}$  with or without data calibration, where the dashed line denotes the event of ending infection with  $\mathcal{R}_t^{14} = 1$ .

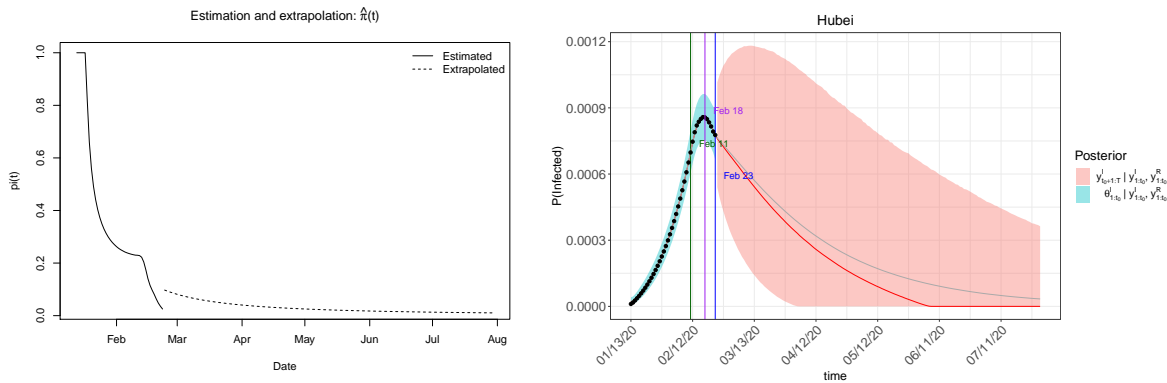


Figure 2: The estimated and predicted mean (grey) and median (red) infection prevalence (right) with transmission rate modifiers  $\hat{\pi}(t)$  (left) fitted by eSIR model using data up to February 18, 2020.

our epidemiological forecast SIR model (Wang et al., 2020). We set  $D = 14$  (or two weeks or the infectious duration that sums the average incubation period and hospitalization time) and computed  $\mathcal{R}_t^{14}$  with or without data calibration used to deal with an abrupt jump of infected cases on February 12, which is mostly due to the issue of under-reporting before February 12. Figure 1 indicates that the estimated  $\mathcal{R}_t^{14}$  appears much smoother (the right panel) with no artificial bumps for the one obtained with the correction for under-reporting than the other obtained with no correction for the under-reporting (the left panel). Note that the reported turning date in Hubei is February 18, which is roughly coincides with the estimated  $\mathcal{R}_t^{14} = 1$ ; see the crossing point of the dashed line and the estimated  $\mathcal{R}_t^{14}$  in Figure 1.

Given that the proportion of cumulative infected cases is very small in comparison to the size of population, one may assume that  $s(t) = S_t/N \approx 1$ . In this case, we may rewrite the transmission rate  $\beta_t$  as a multiplicative form  $\beta_t = \beta_0\pi(t)$ , in which  $\beta_0$  is the transmission rate and  $\pi(t)$  is termed as a transmission rate modifier (Wang et al., 2020). This modifier may be thought of as a consequence of social distancing, self-quarantine, and other preventive interventions issued



during the period of the epidemic in Hubei.

It follows that we obtain an estimate of the transmission rate modifier  $\pi(t)$  as follows:  $\hat{\pi}(t) = \hat{\beta}_t/\beta_0$ . The estimated  $\hat{\pi}(t)$  and its extrapolation are shown in the left panel of Figure 2. we ran the eSIR model (Wang et al., 2020) and obtained the projected infection dynamics in the right panel of Figure 2. Based on results from the eSIR model, we obtained the estimated parameters and their 95% confidence intervals, respectively: the transmission rate  $\hat{\beta}_0 = 0.123$  (CI: [0.0422, 0.256]), the removal rate  $\hat{\gamma} = 0.0257$  (CI: [0.0144, 0.0389]), and the basic reproduction number  $\mathcal{R}_0 = 4.71$  (CI: [2.20, 8.60]) adjusting the time-varying transmission modifier due to various preventive measures implemented in Hubei Province.

In summary, the authors developed a new framework from which one can estimate a time-varying transmission rate and consequently a time-varying effective reproduction number. Their method can provide a way to derive a transmission rate modifier useful in our eSIR model (Wang et al., 2020).

## References

- Wang L, Zhou Y, He J, Zhu B, Wang F, Tang L, et al. (2020). An epidemiological forecast model and software assessing interventions on COVID-19 epidemic in China (with discussion). *Journal of Data Science*, 18(3): 409–454.

## Rejoinder: “Tracking Reproductivity of COVID-19 Epidemic with Varying Coefficient SIR Model”

HAOXUAN SUN<sup>1</sup>, YUMOU QIU<sup>2</sup>, HAN YAN<sup>3</sup>, YAXUAN HUANG<sup>4</sup>, YURU ZHU<sup>5</sup>, JIA GU<sup>5</sup>, AND  
SONG XI CHEN<sup>\*5,6</sup>

<sup>1</sup>*Center for Data Science, Peking University, Beijing, China*

<sup>2</sup>*Department of Statistics, Iowa State University, Ames, Iowa, USA*

<sup>3</sup>*School of Mathematical Sciences, Sichuan University, Chengdu, Sichuan, China*

<sup>4</sup>*Yuanpei College, Peking University, Beijing, China*

<sup>5</sup>*Center for Statistical Science, Peking University, Beijing, China*

<sup>6</sup>*Guanghua School of Management, Peking University, Beijing, China*

We thank the valuable comments from the research team (Wang et al.) from University of Michigan led by Professor Peter Song, the team (Jiang et al.) from Sun Yat-Sen University and University of Science and Technology of China led by Professor Xueqin Wang, and Professor Lu Tang from University of Pittsburgh. These comments broaden the scope of our work.

**Contact Rate Estimation** We are very pleased to see the extension by Wang et al. for estimating the contact rate function  $\pi(t)$  from the estimated time-varying infection rate function  $\hat{\beta}(t)$ . Indeed, the contact rate function is directly reflective to the effects of the early COVID-19 control measures in China which were largely designed to reduce human contacts.

**Euler and Higher Order Approximations to ODEs** Professor Tang’s comments on more accurate discretization of the ODEs are valuable. While the estimator for  $\beta(t)$  is directly motivated by the ODEs of the vSIR model, it can be also formulated under the discrete Poisson-vSIR model under a fixed  $\Delta t$  (which is daily in our study) as indicated before (1) of the paper. Hence, from the prospect of the statistical Poisson-vSIR model, there is no need to conduct the high order correction to the ODEs, although for practical performance such implementation can reduce the bias of the estimation.

**Parametric Form of Infection Rate** The comment by Professor Tang regarding the estimation of the infection rate  $\beta(t)$  as parametric was perhaps based on the reciprocal model (5) proposed for forecasting purposes (to project the future course of the infection rates). Our in-sample estimator of  $\beta(t)$  is the nonparametric kernel estimator.

**Gradual Effects of Abrupt Lockdowns** Regarding Professor Tang’s comments on the smoothed  $\beta(t)$  estimates despite many provinces took sudden preventive actions, the epidemic system from contacts to infections then to diagnosis is a convoluted delayed process. The sudden application of the lock-downs is smoothed by the delayed effects of infections, and the time gap from the onset of the disease to diagnosis, which would make the observed  $\beta(t)$  smooth without sudden drops.

**Moving Average Filters and the Big Revision of Statistics** The moving average filter is applied to remove measurement errors. However, the moving average filter is incapable to deal

---

\*The corresponding authors are Song Xi Chen (csx@gsm.pku.edu.cn) and Yumou Qiu (yumouqiu@iastate.edu).

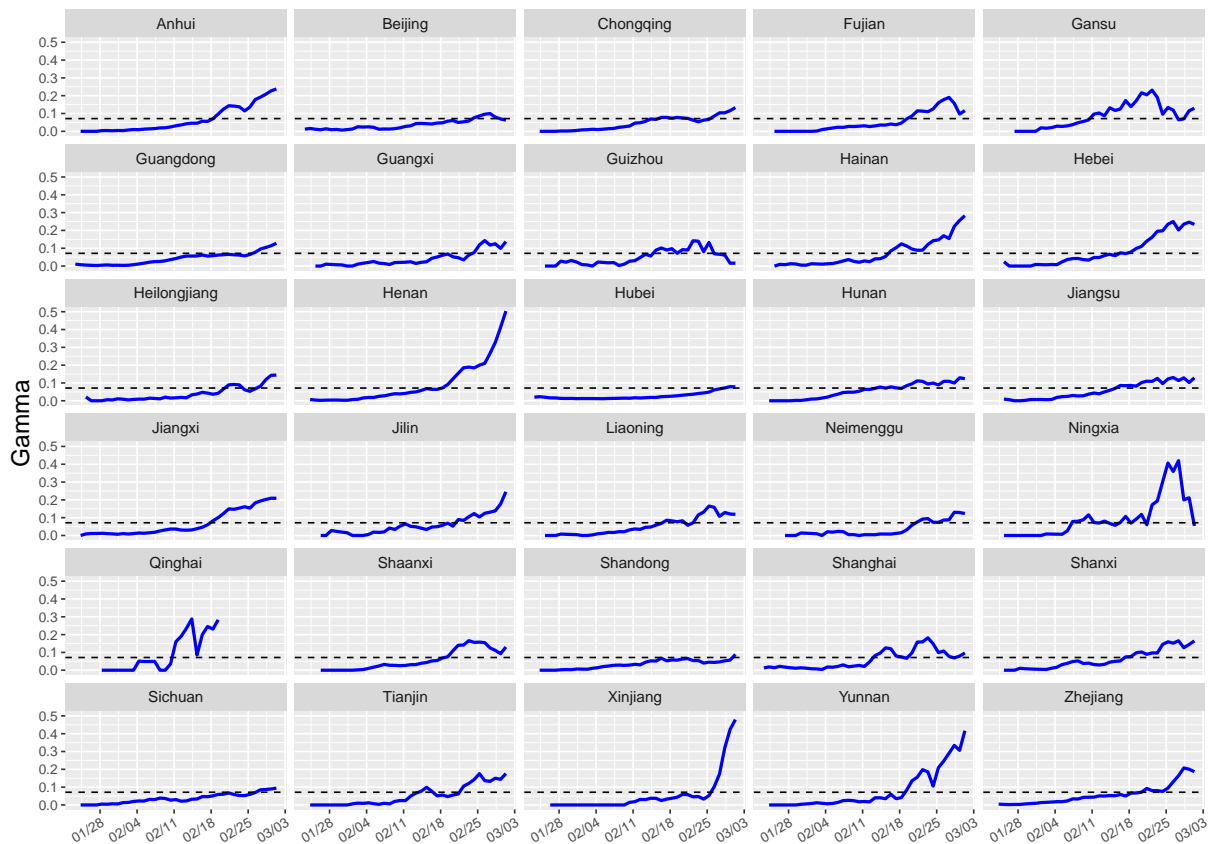


Figure 1: The estimated  $\hat{\gamma}(t)$  from the varying coefficient SIR model (solid) and  $1/14$  (dashed) for the data to Mar 2nd 2020 for 30 provinces.

with the big revision of statistics on February 12th, 2020 in Wuhan. As shown in Section 5, we had applied a one-side linear filter that re-distributes the spikes in the Hubei cities to the previous 7 days with decreasing weights ranging from  $7/28$  to  $1/28$ .

**Roles of Infectious Duration  $D$  on  $R_t$**  Professor Tang's question on  $D$ 's role on  $R_t$  and its variation in  $D$  is timely and important. A major challenge in the estimation of  $R_t$  is the volatility of  $\gamma(t)$  in the early stage of the epidemic. We adopt the formulation of  $R_t^D = \beta(t)D$ . Rather than estimating  $D$ , we assign three sets of  $D$  values in our analysis,  $D = 7, 10.5$  and  $14$ . The narrow confidence intervals for  $R_t$  observed by Wang et al. are likely due to the use of the fixed  $D$ . A wider confidence intervals would appear by adopting the  $D_t = \beta(t)/\gamma(t)$  version via the parametric bootstrap. We have updated Figure 3 (See Figure 1) on the removal rates  $\gamma(t)$  to March 2nd 2020, which shows a general increasing trend in the infection rates.

Regarding the possible over-estimation of  $R_t^{14}$  in Hubei cities other than Wuhan raised by Jiang et al., we would say that the idea of using three  $D$  values to provide a range of  $R_t^D$  values that  $R_t^{14}$  serves as the worse case scenario and hence may over-estimate. In this case, one may use  $R_t^{10.5}$ . Another matter that is in play in the early stage of epidemic is the high volatility in the infection rate estimation despite of adopting a fixed  $D$  version, which may explain the high volatility and the diverse ranges of  $R_t$  estimates in the literature.

**Infection in the Incubation Period and Population Mobility.** As rightly pointed out by Jiang et al., the vSIR framework, which consists of only three compartments, does not accommodate infection in the incubation period, neither the four compartments SEIR model. In a following up work (Gu et al., 2020) we have proposed an extended SEIR model, called the vSEIdR model, which allows infections in both the exposed and infection states. Population mobility or migrations are also not considered, which can be accommodated by adding an extra term to the susceptible population counts  $S(t)$  to reflect immigration or emigration respectively, and a certain percentage of the immigrants (emigrants) are infected patients (imported or exported cases). Doing so would require population mobility data. We hope to conduct such research in a future project.

## References

Gu J, Yan H, Huang Y, Zhu Y, Sun H, Zhang X, et al. (2020). Better strategies for containing COVID-19 epidemics: A study of 25 countries via an extended varying coefficient SEIR model. MedRxiv preprint: <https://doi.org/10.1101/2020.04.27.20081232>.