BETS: The dangers of selection bias in early analyses of the coronavirus disease (COVID-19) pandemic

Qingyuan Zhao

Statistical Laboratory, University of Cambridge

June 15, 2021 @ Queletet Seminar, Ghent University

Paper DOI: 10.1214/20-A0AS1401 Slides: http://www.statslab.cam.ac.uk/~qz280/.

Collaborators



Nianqiao (Phyllis) Ju PhD student at Harvard



Sergio Bacallado Stats Lab, Cambridge



Rajen Shah Stats Lab, Cambridge

And many thanks to...

Cindy Chen, Yang Chen, Yunjin Choi, Hera He, Michael Levy, Marc Lipsitch, James Robins, Andrew Rosenfeld, Dylan Small, Yachong Yang, Zilu Zhou, and many other who have provided helpful suggestions.

• On December 30, 2019, I saw rumours amout a new SARS-like virus on the internet.

- On December 30, 2019, I saw rumours amout a new SARS-like virus on the internet.
- On January 29, 2020, I heard from my parents that a close relative was just diagnosed with "viral pneumonia". This prompted me to start looking into the data available at the time.

- On December 30, 2019, I saw rumours amout a new SARS-like virus on the internet.
- On January 29, 2020, I heard from my parents that a close relative was just diagnosed with "viral pneumonia". This prompted me to start looking into the data available at the time.
- However, epidemiological data from Wuhan are very unreliable!

- On December 30, 2019, I saw rumours amout a new SARS-like virus on the internet.
- On January 29, 2020, I heard from my parents that a close relative was just diagnosed with "viral pneumonia". This prompted me to start looking into the data available at the time.
- However, epidemiological data from Wuhan are very unreliable!

Some anecdotal evidence

- **Inadequate testing:** The relative of mine could not get a RT-PCR test till mid-February, when she was already recovering.
- False negative test: Her first test was negative. A few days later she was tested again and the result came back positive.
- Insufficient contact tracing: Her husband who also showed COVID symptoms quickly recovered and was never tested.



A change of diagnostic criterion on February 12 led to a huge spike of cases.



A change of diagnostic criterion on February 12 led to a huge spike of cases.

Solution: Using cases "exported" from Wuhan



A change of diagnostic criterion on February 12 led to a huge spike of cases.

Solution: Using cases "exported" from Wuhan

This has two benefits:

- Testing and contact tracing were intensive in other locations.
- **②** Detailed case reports (instead of mere case counts) are often available.



A change of diagnostic criterion on February 12 led to a huge spike of cases.

Solution: Using cases "exported" from Wuhan

This has two benefits:

- Testing and contact tracing were intensive in other locations.
- **②** Detailed case reports (instead of mere case counts) are often available.

This design was first used by Neil Ferguson's team in Imperial College, who estimated on January 17 that there might be already over 1,700 cases in Wuhan.

Our first analysis

CSH Spring BMJ Yale	HOME ABO	DUT SUBMIT ALERT:
	Search	
	Comment on this paper	O Previous
of the early 2019-nCoV outbrea	ak using	Posted February 09, 2020.
	of the early 2019-nCoV outbrea	Commercian dis paper and the early 2019-nCoV outbreak using

Our first analysis

medRxiv	CSH String BMJ Yale	HOME AB	OUT SUBMIT ALERT:
Analysis of the epidemic growth	of the early 2019-nCoV outbr	Comment on this paper	Previous Posted February 09, 2020.
Digging Characteristics (Control Cases) Digging Characteristic	11		Download PDF

Methods: We obtained information on the 46 coronavirus cases who traveled from Wuhan before January 23 and have been subsequently confirmed in Hong Kong, Japan, Korea, Macau, Singapore, and Taiwan as of February 5, 2020. Most cases have detailed travel history and disease progress. Compared to previous analyses, an important distinction is that we used this

Our first analysis

medRviv	(CSH) Spring BMJ Yale	HOME ABO	DUT SUBMIT ALERT
THE PREPRINT SERVER FOR HEALTH SCIENCES	Laboratory #	Search	
		Comment on this paper	Previous
Analysis of the epidemic growth o internationally confirmed cases	of the early 2019-nCoV outbrea	ak using	Posted February 09, 2020.
© Qingyuan Zhao, Yang Chen, Dylan S Small doi: https://doi.org/10.1101/2020.02.06.20020941			Download PDF

Methods: We obtained information on the 46 coronavirus cases who traveled from Wuhan before January 23 and have been subsequently confirmed in Hong Kong, Japan, Korea, Macau, Singapore, and Taiwan as of February 5, 2020. Most cases have detailed travel history and disease progress. Compared to previous analyses, an important distinction is that we used this

Results: We found that our model provides good fit to the distribution of the infection time. Assuming the travel rate to the selected countries and regions is constant over the study period, we found that the epidemic was doubling in size every 2.9 days (95% credible interval [CrI], 2 days—4.1 days). Using previously reported serial interval for 2019-nCoV, the estimated basic

A puzzling comparison

THE LANCET



A puzzling comparison

THE LANCET

	ARTICLES VOLUME 395, ISSUE 10225, P689-697, FEBRUARY 29, 2020
	Nowcasting and forecasting the potential domestic and international spread of the 2019-nCoV outbreak originating in Wuhan, China: a
<	modelling study
	Prof Joseph T Wu, PhD 🔌 * 🖾 + Kathy Leung, PhD * - Prof Gabriel M Leung, MD + Show footnotes Published: January 31, 2020 DOI: https://doi.org/10.1016/S0140-6736(20)30260-9 + (1) Check for updates

Methods We used data from Dec 31, 2019, to Jan 28, 2020, on the number of cases exported from Wuhan internationally (known days of symptom onset from Dec 25, 2019, to Jan 19, 2020) to infer the number of infections in Wuhan from Dec 1, 2019, to Jan 25, 2020, Cases exported domestically were then estimated. We forecasted the national and global spread of 2019-nCoV, accounting for the effect of the metropolitan-wide quarantine of Wuhan

A puzzling comparison

THE LANCET



Methods We used data from Dec 31, 2019, to Jan 28, 2020, on the number of cases exported from Wuhan internationally (known days of symptom onset from Dec 25, 2019, to Jan 19, 2020) to infer the number of infections in Wuhan from Dec 1, 2019, to Jan 25, 2020, Cases exported domestically were then estimated. We forecasted the national and global spread of 2019-nCoV, accounting for the effect of the metropolitan-wide quarantine of Wuhan

Findings In our baseline scenario, we estimated that the basic reproductive number for 2019-nCoV was 2-68 (95% Crl 2-47-2-86) and that 75 815 individuals (95% Crl 37304-130330) have been infected in Wuhan as of Jan 25, 2020. The epidemic doubling time was 6-4 days (95% Crl 5-8-7-1). We estimated that in the baseline scenario, Chongqing, Beijing, Shanghai, Guangzhou, and Shenzhen had imported 461 (95% Crl 227-805),

Which one is correct?



In countries most hard hit by COVID-19, the total cases and deaths grew about 100 times in the first 20 days (doubling time: $20/\log_2(100) = 3.01$ days).

How can the results be so different?

Spoilers...

Similar data and model were used in these two studies, with one crucial difference:

The Lancet study did not take into account the travel ban.



Business Markets World Politics TV More



Wuhan lockdown 'unprecedented', shows commitment to contain virus: WHO representative in China

Selection bias in COVID-19 analyses

Types of selection bias

- (i) Under-ascertainment.
- (ii) Non-random sample selection.
- (iii) Travel ban.
- (iv) Epidemic growth.
- (v) Right-truncation.

Selection bias in COVID-19 analyses

Types of selection bias

- (i) Under-ascertainment.
- (ii) Non-random sample selection.
- (iii) Travel ban.
- (iv) Epidemic growth.
- (v) Right-truncation.

Keys to avoid the selection bias

- Carefully design the study and adhere to the sample inclusion criterion.
- Start from a generative model and derive likelihood functions that adjust for sample selection.

Rest of the talk



2 Model

3 Why some early analyses were severely biased?

4 Conclusions

Outline

Dataset

2 Model

- Full data BETS model
- Sample selection
- Results for the parametric model

3 Why some early analyses were severely biased?

- Epidemic growth
- Incubation period

Conclusions

Data collection



- 14 locations where the local health agencies published full case reports.
- 1,460 COVID-19 cases that were confirmed by February 29 for locations in mainland China (February 15 for international locations).

Qingyuan Zhao (Stats Lab, Cambridge)

BETS on COVID-19

Overview of the dataset

Column name	Description	Example	Summary statistics
Case Residence	Unique identifier for each case Nationality or residence of the case	HongKong-05 Wuhan	1460 in total 21.5% reside in Wuhan
Gender	Gender	Male /Female	52.1%/47.7% (0.2% NA)
Age	Age	63	Mean=45.6, IQR=[34, 57]
Known Contact	Known epidemiological contact?	Yes /No	84.7%/15.3%
Cluster	Relationship with other cases	Husband of	32.1% known
		HongKong-04	
Outside	Transmitted outside Wuhan?	Yes/ Likely /No	58.5%/7.7%/33.8%
Begin Wuhan	Begin of stay in Wuhan (B)	30-Nov ⁴	
End Wuhan	End of stay in Wuhan (E)	22-Jan	
Exposure	Period of exposure	1-Dec to 22-Jan	58.9% known period/date 8.2% known date
Arrived	Final arrival date at the location where confirmed a COVID-19 case	22-Jan	40.6% did not travel
Symptom	Date of symptom onset (S)	23-Jan	9.0% NA
Initial	Date of first medical visit	23-Jan	6.5% NA
Confirmed	Date confirmed	24-Jan	

Discerning Wuhan-exported cases

We obtained 378 cases exported from Wuhan that satisfy the following criteria:

- The case had stayed in Wuhan before January 23.
- The case had no recorded contact with other confirmed cases, or had the earliest symptom onset in their (family) cluster, or showed symptoms before they left Wuhan.
- The case did not have missing symptom onset.
- The case arrived at the location where they were diagnosed before January 24.

The principle is to only include cases as Wuhan-exported that pass a **"beyond a reasonable doubt"** test.

Outline

1) Dataset

2 Model

- Full data BETS model
- Sample selection
- Results for the parametric model

3 Why some early analyses were severely biased?

- Epidemic growth
- Incubation period

Conclusions

Outline

Dataset

2 Model

Full data BETS model

- Sample selection
- Results for the parametric model

3 Why some early analyses were severely biased?

- Epidemic growth
- Incubation period

Conclusions

A generative model

Four crucial epidemiological events

- B: Beginning of stay in Wuhan;
- E: End of stay in Wuhan;
- T: Time of transmission (unobserved);
- *S*: Time of symptom onset.

A generative model

Four crucial epidemiological events

- B: Beginning of stay in Wuhan;
- E: End of stay in Wuhan;
- T: Time of transmission (unobserved);
- S: Time of symptom onset.

Below we will:

- Define the support \mathcal{P} of (B, E, T, S) for the Wuhan-exposed population;
- Construct a generative model for (B, E, T, S);
- Define the sample selection set \mathcal{D} corresponds to Wuhan-exported cases;
- Derive likelihood functions to adjust for the sample selection.

Wuhan-exposed population \mathcal{P}

Intuitively, $\mathcal{P} = \text{All people who stayed in Wuhan between 12am December 1, 2019}$ (time 0) and 12am January 24, 2020 (time *L*, the lockdown).

Intuitively, $\mathcal{P} = All$ people who stayed in Wuhan between 12am December 1, 2019 (time 0) and 12am January 24, 2020 (time *L*, the lockdown).

Conventions

• B = 0: Started their stay in Wuhan before time 0.

Intuitively, $\mathcal{P} = \text{All people who stayed in Wuhan between 12am December 1, 2019}$ (time 0) and 12am January 24, 2020 (time *L*, the lockdown).

Conventions

- B = 0: Started their stay in Wuhan before time 0.
- *E* = ∞: Did not arrive in the 14 locations we are considering before time *L*. (We do not differentiate between people who stayed in Wuhan or went to a different location).

Intuitively, $\mathcal{P} = All$ people who stayed in Wuhan between 12am December 1, 2019 (time 0) and 12am January 24, 2020 (time *L*, the lockdown).

Conventions

- B = 0: Started their stay in Wuhan before time 0.
- *E* = ∞: Did not arrive in the 14 locations we are considering before time *L*. (We do not differentiate between people who stayed in Wuhan or went to a different location).
- *T* = ∞: Were not infected during their stay in Wuhan. (We do not differentiate between infection outside Wuhan and never infected.)

Intuitively, $\mathcal{P} = All$ people who stayed in Wuhan between 12am December 1, 2019 (time 0) and 12am January 24, 2020 (time *L*, the lockdown).

Conventions

- B = 0: Started their stay in Wuhan before time 0.
- *E* = ∞: Did not arrive in the 14 locations we are considering before time *L*. (We do not differentiate between people who stayed in Wuhan or went to a different location).
- *T* = ∞: Were not infected during their stay in Wuhan. (We do not differentiate between infection outside Wuhan and never infected.)
- $S = \infty$: Did not show symptoms of COVID-19 (never infected or asymptomatic).

Wuhan-exposed population \mathcal{P}

Intuitively, $\mathcal{P} = All$ people who stayed in Wuhan between 12am December 1, 2019 (time 0) and 12am January 24, 2020 (time *L*, the lockdown).

Conventions

- B = 0: Started their stay in Wuhan before time 0.
- *E* = ∞: Did not arrive in the 14 locations we are considering before time *L*. (We do not differentiate between people who stayed in Wuhan or went to a different location).
- *T* = ∞: Were not infected during their stay in Wuhan. (We do not differentiate between infection outside Wuhan and never infected.)
- $S = \infty$: Did not show symptoms of COVID-19 (never infected or asymptomatic).

Under these conventions.

$$\mathcal{P} = \Big\{(b,e,t,s) \mid b \in [0,L], e \in [b,L] \cup \{\infty\}, t \in [b,e] \cup \{\infty\}, s \in [t,\infty]\Big\}.$$

A generative BETS model

$$f(b, e, t, s) = \underbrace{f_B(b) \cdot f_E(e \mid b)}_{\text{travel}} \cdot \underbrace{f_T(t \mid b, e)}_{\text{disease transmission}} \cdot \underbrace{f_S(s \mid b, e, t)}_{\text{disease progression}}.$$
A generative BETS model

$$f(b, e, t, s) = \underbrace{f_B(b) \cdot f_E(e \mid b)}_{\text{travel}} \cdot \underbrace{f_T(t \mid b, e)}_{\text{disease transmission}} \cdot \underbrace{f_S(s \mid b, e, t)}_{\text{disease progression}}.$$

To allow extrapolation from Wuhan-exported sample to Wuhan-exposed population, the BETS model makes two basic assumptions

Assumption 1: Disease transmission independent of travel

$$f_T(t \mid b, e) = egin{cases} g(t), & ext{if } b < t < e, \ 1 - \int_b^e g(x) \, dx, & ext{if } t = \infty. \end{cases}$$

Here $g(\cdot)$ models the **epidemic growth** in Wuhan before the lockdown.

A generative BETS model

$$f(b, e, t, s) = \underbrace{f_B(b) \cdot f_E(e \mid b)}_{\text{travel}} \cdot \underbrace{f_T(t \mid b, e)}_{\text{disease transmission}} \cdot \underbrace{f_S(s \mid b, e, t)}_{\text{disease progression}} .$$

To allow extrapolation from Wuhan-exported sample to Wuhan-exposed population, the BETS model makes two basic assumptions

Assumption 1: Disease transmission independent of travel

$$f_T(t \mid b, e) = egin{cases} g(t), & ext{if } b < t < e \ 1 - \int_b^e g(x) \, dx, & ext{if } t = \infty. \end{cases}$$

Here $g(\cdot)$ models the **epidemic growth** in Wuhan before the lockdown.

Assumption 2: Disease progression independent of travel

$$f_{S}(s \mid b, e, t) = \begin{cases} \nu \cdot h(s-t), & \text{if } t < s < \infty, \\ 1-\nu, & \text{if } s = \infty. \end{cases}$$

Here $h(\cdot)$ is the density of the **incubation period** S - T (for symptomatic cases).

Parametric assumptions

To ease the interpretation and simply the likelihood functions, we assume

Assumption 3: Exponential growth

$$g(t) = g_{\kappa,r}(t) \stackrel{\Delta}{=} \kappa \cdot \exp(rt), \ t \leq L,$$

Assumption 4: Gamma-distributed incubation period

$$h(s-t) = h_{\alpha,\beta}(s-t) \stackrel{\Delta}{=} rac{eta^{lpha}}{\Gamma(lpha)}(s-t)^{lpha-1} \exp\{-eta(s-t)\}.$$

Parametric assumptions

To ease the interpretation and simply the likelihood functions, we assume

Assumption 3: Exponential growth

$$g(t) = g_{\kappa,r}(t) \stackrel{\Delta}{=} \kappa \cdot \exp(rt), \ t \leq L,$$

Assumption 4: Gamma-distributed incubation period

$$h(s-t) = h_{\alpha,\beta}(s-t) \stackrel{\Delta}{=} rac{eta^{lpha}}{\Gamma(lpha)}(s-t)^{lpha-1} \exp\{-eta(s-t)\}.$$

- The nuisance parameters ν (proportion of symptomatic cases) and κ (baseline transmission) will be canceled in the likelihood function.
- Assumptions 3 & 4 are relaxed in a Bayesian nonparametric analysis (can be found in the paper).

Outline

Dataset

2 Model

- Full data BETS model
- Sample selection
- Results for the parametric model

3 Why some early analyses were severely biased?

- Epidemic growth
- Incubation period

4 Conclusions

The event of observing Wuhan-exported cases can be written as

$$\mathcal{D} = \{ (b, e, t, s) \in \mathcal{P} \mid b \le t \le e \le L, t \le s < \infty \}.$$

The event of observing Wuhan-exported cases can be written as

$$\mathcal{D} = \{ (b, e, t, s) \in \mathcal{P} \mid b \le t \le e \le L, t \le s < \infty \}.$$

This makes three further restrictions on \mathcal{P} :

The event of observing Wuhan-exported cases can be written as

$$\mathcal{D} = \{(b, e, t, s) \in \mathcal{P} \mid b \leq t \leq e \leq L, t \leq s < \infty\}.$$

This makes three further restrictions on \mathcal{P} :

B ≤ T ≤ E, because we only use cases who contracted the virus during their stay in Wuhan;

The event of observing Wuhan-exported cases can be written as

$$\mathcal{D} = \{(b, e, t, s) \in \mathcal{P} \mid b \leq t \leq e \leq L, t \leq s < \infty\}.$$

This makes three further restrictions on \mathcal{P} :

- B ≤ T ≤ E, because we only use cases who contracted the virus during their stay in Wuhan;
- *E* ≤ *L*, because the case can only be observed if they left Wuhan before the travel ban;

The event of observing Wuhan-exported cases can be written as

$$\mathcal{D} = \{(b, e, t, s) \in \mathcal{P} \mid b \leq t \leq e \leq L, t \leq s < \infty\}.$$

This makes three further restrictions on \mathcal{P} :

- B ≤ T ≤ E, because we only use cases who contracted the virus during their stay in Wuhan;
- *E* ≤ *L*, because the case can only be observed if they left Wuhan before the travel ban;
- § $S < \infty$, because we only consider COVID-19 cases who showed symptoms.

For a moment, let's pretend the time of transmission T is observed.

For a moment, let's pretend the time of transmission T is observed.



X Sample from \mathcal{P}

For a moment, let's pretend the time of transmission T is observed.

$$\prod_{i=1}^n f(B_i, E_i, T_i, S_i)$$

✓ Sample from D (Unconditional likelihood)

$$\prod_{i=1}^{n} f(B_{i}, E_{i}, T_{i}, S_{i} \mid \mathcal{D}), \text{ where } f(b, e, t, s \mid \mathcal{D}) \triangleq \frac{f(b, e, t, s) \cdot 1_{\{(b, e, t, s) \in \mathcal{D}\}}}{\mathbb{P}((B, E, T, S) \in \mathcal{D})}$$

X Sample from \mathcal{P}

For a moment, let's pretend the time of transmission T is observed.

$$\prod_{i=1}^n f(B_i, E_i, T_i, S_i)$$

✓ Sample from D (Unconditional likelihood)

$$\prod_{i=1}^{n} f(B_i, E_i, T_i, S_i \mid \mathcal{D}), \text{ where } f(b, e, t, s \mid \mathcal{D}) \triangleq \frac{f(b, e, t, s) \cdot 1_{\{(b, e, t, s) \in \mathcal{D}\}}}{\mathbb{P}((B, E, T, S) \in \mathcal{D})}$$

✓ Sample from D (Conditional likelihood)

$$\prod_{i=1}^{n} f(T_i, S_i \mid B_i, E_i, \mathcal{D}), \text{ where } f(t, s \mid b, e, \mathcal{D}) \triangleq \frac{f(t, s \mid B = b, E = e) \cdot 1_{\{(b, e, t, s) \in \mathcal{D}\}}}{\mathbb{P}((B, E, T, S) \in \mathcal{D} \mid B = b, E = e)}.$$

In reality, the time of transmission T is unobserved. We can either treat T as a latent variable and use e.g. an EM algorithm, or use the **integrated likelihood**:

In reality, the time of transmission T is unobserved. We can either treat T as a latent variable and use e.g. an EM algorithm, or use the **integrated likelihood**:

Unconditional likelihood

$$L_{\text{uncond}}(\theta) = \prod_{i=1}^{n} \int f(B_i, E_i, t, S_i \mid D) dt,$$

where $\theta = (f_B(\cdot), f_E(\cdot | \cdot), g(\cdot), h(\cdot)).$

In reality, the time of transmission T is unobserved. We can either treat T as a latent variable and use e.g. an EM algorithm, or use the **integrated likelihood**:

Unconditional likelihood

$$L_{uncond}(\theta) = \prod_{i=1}^{n} \int f(B_i, E_i, t, S_i \mid D) dt,$$

where $\theta = (f_B(\cdot), f_E(\cdot | \cdot), g(\cdot), h(\cdot)).$

Conditional likelihood

$$L_{\text{cond}}(\theta) = \prod_{i=1}^{n} \int f(t, S_i \mid B_i, E_i, \mathcal{D}) dt,$$

where $\theta = (g(\cdot), h(\cdot))$.

In reality, the time of transmission T is unobserved. We can either treat T as a latent variable and use e.g. an EM algorithm, or use the **integrated likelihood**:

Unconditional likelihood

$$L_{uncond}(\theta) = \prod_{i=1}^{n} \int f(B_i, E_i, t, S_i \mid D) dt,$$

where $\theta = (f_B(\cdot), f_E(\cdot | \cdot), g(\cdot), h(\cdot)).$

Conditional likelihood

$$L_{\text{cond}}(\theta) = \prod_{i=1}^{n} \int f(t, S_i \mid B_i, E_i, \mathcal{D}) dt,$$

where $\theta = (g(\cdot), h(\cdot))$.

The conditional likelihood is less efficient because it does not use information in $f(b, e \mid D)$; but it is robust to misspecifying the travel models $f_B(\cdot), f_E(\cdot \mid \cdot)$.

Conditional likelihood function

Proposition

Under Assumptions 1-4,

$$\begin{split} & L_{\text{cond}}(r,\alpha,\beta) = \\ & \left\{ r^n \Big(\frac{\beta}{\beta+r} \Big)^{n\alpha} \cdot \prod_{i=1}^n \frac{\exp(rS_i) \big[H_{\alpha,\beta+r}(S_i - B_i) - H_{\alpha,\beta+r}((S_i - E_i)_+) \big]}{\exp(rE_i) - \exp(rB_i)}, & \text{for } r > 0, \\ & \prod_{i=1}^n \frac{H_{\alpha,\beta}(S_i - B_i) - H_{\alpha,\beta}((S_i - E_i)_+)}{E_i - B_i}, & \text{for } r = 0, \end{split} \right. \end{split}$$

where $H_{\alpha,\beta}(\cdot)$ is the CDF of Gamma (α,β) and $(\cdot)_+ = \max(\cdot,0)$ is the positive part function.

Conditional likelihood function

Proposition

Under Assumptions 1-4,

$$\begin{split} & \mathcal{L}_{\text{cond}}(r,\alpha,\beta) = \\ & \begin{cases} r^n \Big(\frac{\beta}{\beta+r}\Big)^{n\alpha} \cdot \prod_{i=1}^n \frac{\exp(rS_i) \big[H_{\alpha,\beta+r}(S_i - B_i) - H_{\alpha,\beta+r}((S_i - E_i)_+) \big]}{\exp(rE_i) - \exp(rB_i)}, & \text{for } r > 0, \\ & \prod_{i=1}^n \frac{H_{\alpha,\beta}(S_i - B_i) - H_{\alpha,\beta}((S_i - E_i)_+)}{E_i - B_i}, & \text{for } r = 0, \end{cases} \end{split}$$

where $H_{\alpha,\beta}(\cdot)$ is the CDF of Gamma (α,β) and $(\cdot)_+ = \max(\cdot,0)$ is the positive part function.

- Does not depend on ν (proportion of symptomatic cases) and κ (baseline transmission).
- When r = 0, reduces to the likelihood function in Reich et al. (2009) *Statistics in Medicine*, 28:2769–2784.

Unconditional likelihood function

Assumption 5: Stable travel

- **(**) Beginning of stay *B* follows a uniform distribution given $0 < B \leq L$.
- End of stay E follows a uniform distribution from B to L (with different rates for Wuhan residents and Wuhan visitors).

Unconditional likelihood function

Assumption 5: Stable travel

- **(**) Beginning of stay *B* follows a uniform distribution given $0 < B \leq L$.
- End of stay E follows a uniform distribution from B to L (with different rates for Wuhan residents and Wuhan visitors).

Proposition

Under Assumptions 1-5 and suitable approximations,

$$\begin{split} L_{\text{uncond}}(\rho, r, \alpha, \beta) &\approx r^{2n} \Big(\frac{\beta}{\beta+r}\Big)^{n\alpha} \cdot \prod_{i=1}^{n} \bigg\{ \frac{\mathbb{1}_{\{B_i=0\}} + (\rho/L)\mathbb{1}_{\{B_i>0\}}}{\mathbb{1} + \rho(\mathbb{1} - 2/(rL))} \exp\big\{r(S_i - L)\big\} \\ &\times \big[H_{\alpha,\beta+r}(S_i - B_i) - H_{\alpha,\beta+r}((S_i - E_i)_+)\big]\bigg\}, \end{split}$$

where ρ is a traveling parameter (capturing the different traveling patterns between Wuhan residents and visitors).

Outline

Dataset

2 Model

- Full data BETS model
- Sample selection
- Results for the parametric model

3 Why some early analyses were severely biased?

- Epidemic growth
- Incubation period

Conclusions

Results

Location	Sample size	Doubling time (in days)	Incubation period Median 95% quantile	
Conditional likelihood				
China - Hefei	34	2.1 (1.2–3.7)	4.3 (2.9-6.0)	12.0 (9.1–17.3)
China - Shaanxi	53	1.7 (1.0-2.8)	4.5 (3.1-6.2)	14.6 (11.5–19.8)
China - Shenzhen	129	2.2 (1.7-3.0)	3.5 (2.8–4.3)	11.2 (9.5–13.6)
China - Xinyang	74	2.3 (1.5–3.5)	6.8 (5.4–8.2)	16.4 (13.8–20.1)
China - Other	42	2.0 (1.1–3.4)	5.1 (3.6-6.7)	12.3 (9.8–16.4)
International	46	2.1 (1.4–3.4)	3.8 (2.5–5.3)	10.9 (8.4–15.1)
All locations	378	2.1 (1.8–2.5)	4.5 (4.0-5.0)	13.4 (12.2–14.8)
Unconditional likelihood				
China - Hefei	34	1.8 (1.4–2.4)	4.1 (2.8–5.5)	11.9 (9.0-17.2)
China - Shaan×i	53	2.5 (2.0-3.1)	5.3 (3.9–6.8)	15.0 (12.0-20.0)
China - Shenzhen	129	2.4 (2.1–2.8)	3.6 (2.9-4.3)	11.3 (9.6–13.7)
China - Xinyang	74	2.4 (2.0-2.9)	6.8 (5.6-8.1)	16.4 (13.9–20.2)
China - Other	42	2.1 (1.7–2.8)	5.3 (4.0-6.6)	12.4 (10.0-16.4)
International	46	2.0 (1.6-2.6)	3.7 (2.5–5.0)	10.8 (8.4–15.1)
All locations	378	2.3 (2.1–2.5)	4.6 (4.1–5.1)	13.5 (12.3–14.9)

(Point estimates obtained by MLE. Confidence intervals obtained by inverting LRT.)

Conclusions from the parametric model

- The initial doubling time in Wuhan is between 2 to 2.5 days.
- The median incubation period is around 4 days.
- The 95% quantile of the incubation period is between 11 to 15 days.

Outline

Dataset

2 Model

- Full data BETS model
- Sample selection
- Results for the parametric model

3 Why some early analyses were severely biased?

- Epidemic growth
- Incubation period

Conclusions

Outline

1 Dataset

2 Model

- Full data BETS model
- Sample selection
- Results for the parametric model

3 Why some early analyses were severely biased?

- Epidemic growth
- Incubation period

Conclusions

A puzzling comparison

THE LANCET



Methods We used data from Dec 31, 2019, to Jan 28, 2020, on the number of cases exported from Wuhan internationally (known days of symptom onset from Dec 25, 2019, to Jan 19, 2020) to infer the number of infections in Wuhan from Dec 1, 2019, to Jan 25, 2020, Cases exported domestically were then estimated. We forecasted the national and global spread of 2019-nCoV, accounting for the effect of the metropolitan-wide quarantine of Wuhan

Findings In our baseline scenario, we estimated that the basic reproductive number for 2019-nCoV was 2-68 (95% Crl 2-47-2-86) and that 75 815 individuals (95% Crl 37304-130330) have been infected in Wuhan as of Jan 25, 2020. The epidemic doubling time was 6-4 days (95% Crl 5-8-7-1). We estimated that in the baseline scenario, Chongqing, Beijing, Shanghai, Guangzhou, and Shenzhen had imported 461 (95% Crl 227-805),

What happened?

Wu et al. used a modified SEIR (Susceptible-Exposed-Infectious-Recovered) model to account for traveling. But they **did not consider the travel ban**.

What happened?

Wu et al. used a modified SEIR (Susceptible-Exposed-Infectious-Recovered) model to account for traveling. But they **did not consider the travel ban**.

X Density of S in \mathcal{P}

It is reasonable to assume incidence of symptom onset is growing exponentially in Wuhan-exposed population \mathcal{P} :

```
f(s \mid \mathcal{P}) \underset{\sim}{\propto} \exp(rs), \text{ for } s \leq L.
```

But we are sampling from the Wuhan-exported cases \mathcal{D} .

What happened?

Wu et al. used a modified SEIR (Susceptible-Exposed-Infectious-Recovered) model to account for traveling. But they **did not consider the travel ban**.

X Density of S in \mathcal{P}

It is reasonable to assume incidence of symptom onset is growing exponentially in Wuhan-exposed population \mathcal{P} :

$$f(s \mid \mathcal{P}) \underset{\sim}{\propto} \exp(rs), \text{ for } s \leq L.$$

But we are sampling from the Wuhan-exported cases \mathcal{D} .

✓ Density of *S* in D

Under Assumptions 1-5 and reasonable approximations,

$$f(t \mid \mathcal{D}, B = 0) \underset{\sim}{\propto} \exp(rt) \left(L - t \right) \mathbf{1}_{\{t \leq L\}},$$

We can further derive the theoretical $f_S(s \mid D, B = 0)$; in particular,

$$f_{\mathcal{S}}(s \mid \mathcal{D}, B = 0) \propto \exp(rs) \left(L + \frac{\alpha}{\beta + r} - s \right), \text{ for } s \leq L.$$

Illustration of the selection bias (iii)



- Histogram: Density of the symptom onset of the Wuhan-resident cases;
- Orange curve: Theoretical fit $f_{\mathcal{S}}(s \mid \mathcal{D}, B = 0)$ using MLE of (r, α, β) .
- Blue dashed line: January 23, 2020 (time L).

Outline

Dataset

2 Model

- Full data BETS model
- Sample selection
- Results for the parametric model

3 Why some early analyses were severely biased?

- Epidemic growth
- Incubation period

Conclusions

Bias (iv): Epidemic growth

- Patients were more likely to be infected towards the end of their exposure period.
- **Susceptible studies:** Studies that treat infections as uniformly distributed over the exposure period.
- Direction of bias: Over-estimation of the incubation period.
- Solution: Use the likelihood $L_{cond}(r, \alpha, \beta)$ instead of $L_{cond}(0, \alpha, \beta)$.

Bias (v): Right-truncation

- Cases confirmed after a certain time are excluded from the dataset.
- **Susceptible studies:** Studies that only use cases detected early in an epidemic.
- Direction of bias: Under-estimation of the incubation period.
- Solution: Derive the likelihood with the additional conditioning event $S \leq M$.

Likelihood function adjusted for right-truncation

• Under Assumptions 1 & 2,

$$f_{T,S}(t,s \mid b,e,\mathcal{D},S \leq M) = \frac{g(t)h(s-t)}{\int_b^{\max(e,s)} g(t)H(M-t)\,dt}$$

where $H(\cdot)$ is the CDF of $h(\cdot)$.

 Closed-form expression for L_{cond,trunc}(r, α, β; M) can further be obtained under Assumptions 3 & 4 using integration by parts.

Illustration of the selection bias (iv) and (v)

An experiment

• For each day between January 23 and February 18, obtain the subset of cases confirmed by that day.
Illustration of the selection bias (iv) and (v)

An experiment

- For each day between January 23 and February 18, obtain the subset of cases confirmed by that day.
- Fit the parametric BETS model by using one of the following likelihoods:
 - Adjusted for nothing: L_{cond}(0, α, β) (likelihood function in Reich et al. (2009) used in other studies).
 - **2** Adjusted for growth: $L_{cond}(r, \alpha, \beta)$.
 - **3** Adjusted for growth and right-truncation: $L_{\text{cond,trunc}}(r, \alpha, \beta; M)$.

Illustration of the selection bias (iv) and (v)

An experiment

- For each day between January 23 and February 18, obtain the subset of cases confirmed by that day.
- Fit the parametric BETS model by using one of the following likelihoods:
 - Adjusted for nothing: L_{cond}(0, α, β) (likelihood function in Reich et al. (2009) used in other studies).
 - **2** Adjusted for growth: $L_{cond}(r, \alpha, \beta)$.
 - **3** Adjusted for growth and right-truncation: $L_{\text{cond,trunc}}(r, \alpha, \beta; M)$.
- Obtain point estimates by MLE and CIs by nonparametric Bootstrap.

Illustration of the selection bias (iv) and (v)

An experiment

- For each day between January 23 and February 18, obtain the subset of cases confirmed by that day.
- Fit the parametric BETS model by using one of the following likelihoods:
 - **Adjusted for nothing:** $L_{cond}(0, \alpha, \beta)$ (likelihood function in Reich et al. (2009) used in other studies).
 - **2** Adjusted for growth: $L_{cond}(r, \alpha, \beta)$.
 - **3** Adjusted for growth and right-truncation: $L_{cond,trunc}(r, \alpha, \beta; M)$.
- Obtain point estimates by MLE and CIs by nonparametric Bootstrap.
- Compare with previous studies:
 - Backer, J. A. et al. *Eurosurveillance*, 25(5), 2020. PubMed: 32046819.
 - 2 Lauer, S. A. et al. Annals of Internal Medicine, 2020. PubMed: 32150748.
 - Linton, N. M. et al. Journal of Clinical Medicine, 9(2), 2020. PubMed: 32079150.



Ignore epidemic growth \implies Overestimate incubation period. Ignore right-truncation \implies Underestimate incubation period.

Outline

Dataset

2 Model

- Full data BETS model
- Sample selection
- Results for the parametric model

3 Why some early analyses were severely biased?

- Epidemic growth
- Incubation period

4 Conclusions

Conclusions about COVID-19

- Initial doubling time in Wuhan: 2-2.5 days.
- Median incubation period: about 4 days.
- Proportion of incubation period at least 14 days: about 5%.

Conclusions about COVID-19

- Initial doubling time in Wuhan: 2-2.5 days.
- Median incubation period: about 4 days.
- Proportion of incubation period at least 14 days: about 5%.

Our study has many limitations:

- Reported symptom onset could be inaccurate.
- Some degree of under-ascertainment is perhaps inevitable.
- Discerning Wuhan-exported cases is not black-and-white.
- Assumptions 1 & 2 (independence of travel and disease) could be violated.

Compelling evidence for selection bias in early studies

- (i) Under-ascertainment.
- (ii) Non-random sample selection.
- (iii) Travel ban.
- (iv) Epidemic growth.
- (v) Right-truncation.

Compelling evidence for selection bias in early studies

- (i) Under-ascertainment.
- (ii) Non-random sample selection.
- (iii) Travel ban.
- (iv) Epidemic growth.
- (v) Right-truncation.

Don't make uncalculated BETS

- Carefully design the study and adhere to the sample inclusion criterion.
- Base statistical inference on first principles.

Compelling evidence for selection bias in early studies

- (i) Under-ascertainment.
- (ii) Non-random sample selection.
- (iii) Travel ban.
- (iv) Epidemic growth.
- (v) Right-truncation.

Don't make uncalculated BETS

- Carefully design the study and adhere to the sample inclusion criterion.
- Base statistical inference on first principles.

Final Lesson:

Data Quality + Better Design >> Data Quantity + Better Model