







# Accurately Estimating Unreported Infections using Information Theory

*Jiaming Cui*<sup>1</sup>, Bijaya Adhikari<sup>2</sup>, Arash Haddadan<sup>3</sup>, A S M Ahsan-Ul Haque<sup>3</sup>, Jilles Vreeken<sup>4</sup>, Anil Vullikanti<sup>3</sup>, B. Aditya Prakash<sup>1</sup>

> <sup>1</sup>Georgia Institute of Technology <sup>2</sup>University of Iowa <sup>3</sup>University of Virginia <sup>4</sup>CISPA Helmholtz Center for Information Security



May 3, 2025 @ SDM'25

• One of the most significant challenges in combating against the spread of infectious diseases

- One of the most significant challenges in combating against the spread of infectious diseases
- For example, large number of COVID-19 infections were unreported
  - Lack of testing
  - Asymptomatic infections

- One of the most significant challenges in combating against the spread of infectious diseases
- For example, large number of COVID-19 infections were unreported
  - Lack of testing
  - Asymptomatic infections



Only 23 reported infections in 5 major U.S. cities by March 1, 2020.

- One of the most significant challenges in combating against the spread of infectious diseases
- For example, large number of COVID-19 infections were unreported
  Boston
  - Lack of testing
  - Asymptomatic infections



Only 23 reported infections in 5 major U.S. cities by March 1, 2020.

It was estimated that there were already more than 28000 infections. <sup>5</sup>

- Seattle
- Chicago
- San Francisco
- New York

- One of the most significant challenges in combating against the spread of infectious diseases
- For example, large number of COVID-19 infections were unreported
  Boston
  - Lack of testing
  - Asymptomatic infections

- Seattle
- Chicago
- San Francisco
- New York

Inability in estimating the unreported infections allowed them to drive up disease spread in the U.S. and worldwide.

Only 23 reported infections in 5 It was estimated that there were already major U.S. cities by March 1, 2020. more than 28000 infections. <sup>6</sup>

### **Reported rate**

- Epidemiologists use reported rate (α<sub>reported</sub>) to capture total infections (i.e., both reported and unreported)
- Definition:

 $\alpha_{\text{reported}} = \frac{\text{reported infections}}{\text{total infections}}$ 

### **Reported rate estimation methods**

- Serological studies: Prevalence of antibiotics
  - Accurate, but expensive and have unavoidable delays
- Influenza surveillance systems



- Suffer from ad-hoc corrections between COVID-19 and influenza
- Epidemiological models

### **Reported rate estimation methods**

- Serological studies: Prevalence of antibiotics
  - Accurate, but expensive and have unavoidable delays
- Influenza surveillance systems



- Suffer from ad-hoc corrections between COVID-19 and influenza
- Epidemiological models

- Many epi models have reported rate as a parameter
- Example:



$$\begin{aligned} \frac{dS}{dt} &= -\beta \frac{SI}{N} \\ \frac{dE}{dt} &= \beta \frac{SI}{N} - \gamma E \\ \frac{dI_r}{dt} &= \alpha_{reported} \gamma E - \delta I_r \\ \frac{dI_u}{dt} &= (1 - \alpha_{reported}) \gamma E - \delta I_u \\ \frac{dR}{dt} &= \delta (I_r + I_u) \end{aligned}$$

- Many epi models have reported rate as a parameter
- Example:

States for reported  $(I_r)$  and unreported

cases (I<sub>11</sub>)

$$\begin{aligned} \frac{dS}{dt} &= -\beta \frac{SI}{N} \\ \frac{dE}{dt} &= \beta \frac{SI}{N} - \gamma E \\ \frac{dI_r}{dt} &= \alpha_{reported} \gamma E - \delta I_r \\ \frac{dI_u}{dt} &= (1 - \alpha_{reported}) \gamma E - \delta I_u \\ \frac{dR}{dt} &= \delta (I_r + I_u) \end{aligned}$$

- Many epi models have reported rate as a parameter
- Example:

**States for reported** 

 $(I_r)$  and unreported

cases (I<sub>11</sub>)

$$\begin{aligned} \frac{dS}{dt} &= -\beta \frac{SI}{N} \\ \frac{dE}{dt} &= \beta \frac{SI}{N} - \gamma E \\ \frac{dI_r}{dt} &= \alpha_{reported} \gamma E - \delta I_r \\ \frac{dI_u}{dt} &= (1 - \alpha_{reported}) \gamma E - \delta I_u \\ \frac{dR}{dt} &= \delta (I_r + I_u) \end{aligned}$$

 $\alpha_{reported}$ : reported rate

• Many epi models have reported rate as a parameter



• Calibration: Fit  $I_r$  to reported cases to estimate the unknown parameters (including reported rate)

#### • Parameter tuning is hard

Example: COVID-19 cases in Florida in 2020



#### • Parameter tuning is hard

• Example: COVID-19 cases in Florida in 2020



#### • Parameter tuning is hard

Example: COVID-19 cases in Florida in 2020



#### • Parameter tuning is hard

Example: COVID-19 cases in Florida in 2020



# **Our goal**

• Therefore, our goal is to estimate an **accurate** reported rate

### Intuition

- We already know reported cases D<sub>reported</sub>
- Imagine we are also given accurate values of total cases D
  - Then calibrating the model to  $(D, D_{reported})$  will lead to better fit of  $D_{reported}$
  - We can also learn better parameters including  $\alpha_{reported}$
- Hence, our problem can be stated as finding the *D*\* that fits the *D<sub>reported</sub>* best

# Outline

- Introduction
- Two-part MDL: Sender-receiver framework
- MDLINFER: Information-theory based approach
- Experiment results
- Conclusions

### **Two-part MDL: Sender-receiver framework**

- Two hypothetical actors: Sender *S* and Receiver *R* 
  - Sender *S* wants to send the "DATA" to Receiver *R* using a good "MODEL"
  - To measure the "good fit", use the number of bits to encode DATA: L(DATA|MODEL) + L(MODEL)

### **Two-part MDL: Sender-receiver framework**

- Two hypothetical actors: Sender *S* and Receiver *R* 
  - Sender S wants to send the "DATA" to Receiver R using a good "MODEL"
  - To measure the "good fit", use the number of bits to encode DATA: L(DATA|MODEL) + L(MODEL)



### **Two-part MDL: Sender-receiver framework**

- Two hypothetical actors: Sender *S* and Receiver *R* 
  - Sender S wants to send the "DATA" to Receiver R using a good "MODEL"
  - To measure the "good fit", use the number of bits to encode DATA: L(DATA|MODEL) + L(MODEL)

A good MODEL makes description easier!

Sender *S* searches for the best possible MODEL, which minimizes the overall cost of encoding and transmitting both the MODEL and the DATA given the MODEL.

Describe it L(DATA) = L(MODEL) + L(DATA|MODEL)directly is hard!

# Outline

- Introduction
- Two-part MDL: Sender-receiver framework
- MDLINFER: Information-theory based approach
- Experiment results
- Conclusions

### **MDLINFER: MODEL and DATA**

- Recall that our problem can be stated as finding the D\* that fits the D<sub>reported</sub> best
- We also have

 $D_{\text{reported}} = \alpha_{\text{reported}} \times D$ 

- Hence, we intuitively use
  - $\circ$  *D*<sub>reported</sub> as DATA
  - $\circ$  *D* and  $\alpha_{reported}$  as MODEL

### **MDLINFER: Problem formulation**

Problem formulation

$$D^* = \underset{D}{\operatorname{argmin}} L(D_{reported} | D, \alpha_{reported}) + L(D, \alpha_{reported})$$

• Here,  $L(\cdot)$  denotes the number of bits for encoding



Use MDL Sender-Receiver framework to help find the latent variable  $D^*$  that explains  $D_{reported}$  best. The latent variable  $D^*$  helps find more accurate  $\theta^*$  for the epidemiological model.

### **MDLINFER: Problem formulation**

Problem formulation

$$D^* = \underset{D}{\operatorname{argmin}} L(D_{reported} | D, \alpha_{reported}) + L(D, \alpha_{reported})$$

• Here,  $L(\cdot)$  denotes the number of bits for encoding



# Outline

- Introduction
- Two-part MDL: Sender-receiver framework
- MDLINFER: Information-theory based approach
- Experiment results
- Conclusions

- More accurate estimation of total cases
  - Black: ground-truth total cases by serological studies <sup>[1]</sup>
  - Red for us, blue for current estimation methods



29

- More accurate estimation of total cases
  - Black: ground-truth total cases by serological studies <sup>[1]</sup>
  - Red for us, blue for current estimation methods





- More accurate estimation of total cases
  - Black: ground-truth total cases by serological studies <sup>[1]</sup> Ο
  - Red for us, blue for current estimation methods



[1] https://www.cdc.gov/coronavirus/2019-ncov/cases-updates/commercial-lab-surveys.html

More results for other locations in the paper

### Results

- More accurate estimation of total cases
  - Black: ground-truth total cases by serological studies <sup>[1]</sup>
  - Red for us, blue for current estimation methods



CDC

[1] https://www.cdc.gov/coronavirus/2019-ncov/cases-updates/commercial-lab-surveys.html

- Better fit of reported cases
  - Black for ground-truth reported cases by NY Times <sup>[1]</sup>
  - Red for us, blue for current estimation methods

- Better fit of reported cases
  - Black for ground-truth reported cases by NY Times <sup>[1]</sup>
  - Red for us, blue for current estimation methods



Շի։ <sup>Ջշավի տե</sup> Ծimes [1] https://www.nytimes.com/interactive/2021/us/covid-cases.html

- Better fit of reported cases
  - Black for ground-truth reported cases by NY Times <sup>[1]</sup>
  - Red for us, blue for current estimation methods



More results for other locations in the paper

### Results

- Better fit of reported cases
  - Black for ground-truth reported cases by NY Times <sup>[1]</sup>
  - Red for us, blue for current estimation methods



- Better estimate for symptomatic rate trends
  - Black for symptom trends by Facebook's survey <sup>[1]</sup>
  - Red for us, blue for current estimation methods

37

- Better estimate for symptomatic rate trends
  - Black for symptom trends by Facebook's survey <sup>[1]</sup>
  - Red for us, blue for current estimation methods



38

- Better estimate for symptomatic rate trends
  - Black for symptom trends by Facebook's survey <sup>[1]</sup>
  - Red for us, blue for current estimation methods



[1] Reinhart, Alex, et al. An open repository of real-time COVID-19 indicators. PNAS 2021.

More results for other locations in the paper

### Results

- Better estimate for symptomatic rate trends
  - Black for symptom trends by Facebook's survey <sup>[1]</sup>
  - Red for us, blue for current estimation methods



[1] Reinhart, Alex, et al. An open repository of real-time COVID-19 indicators. PNAS 2021.

# Outline

- Introduction
- Two-part MDL: Sender-receiver framework
- MDLINFER: Information-theory based approach
- Experiment results
- Conclusions

# Conclusions



- We propose MDLINFER, a data-driven method to identify reported rate
- Leverage the information theory-based MDL framework
- Better performance in identifying total infections and forecasting future infections

$$\begin{array}{c} O & O \\ O & O \\ \end{array} = \begin{array}{c} O & O \\ - \begin{array}{c} O & O \\ O \\ \end{array} \\ L(DATA) &= L(MODEL) + L(DATA|MODEL) \end{array}$$

### **Authors**



Jiaming Cui



Bijaya Adhikari





Arash Haddadan A S M Ahsan-Ul Haque



Jilles Vreeken



Anil Vullikanti



B. Aditya Prakash



### Thank you

 Code & papers available at: people.cs.vt.edu/jiamingcui/



Acknowledgements: NSF (Expeditions CCF-1918770 and CCF-1918656, CAREER IIS-2028586, RAPID IIS- 2027862, Medium IIS-1955883, Medium IIS-2106961, IIS-2403240, IIS-1931628, IIS-1955797, IIS- 2027848, IIS-2331315, PIPP CCF-2200269), NIH 2R01GM109718, CDC MInD program 44 U01CK000589, ORNL, Dolby faculty research award, UVA GIDI, Georgia Tech.