

Estimation of COVID-19 prevalence rates using EM algorithm

June 2025

Jungsik Noh


Outline

- COVID-19 data in 2020
- Data science perspective
- Expectation-maximization algorithm

PLOS ONE

RESEARCH ARTICLE

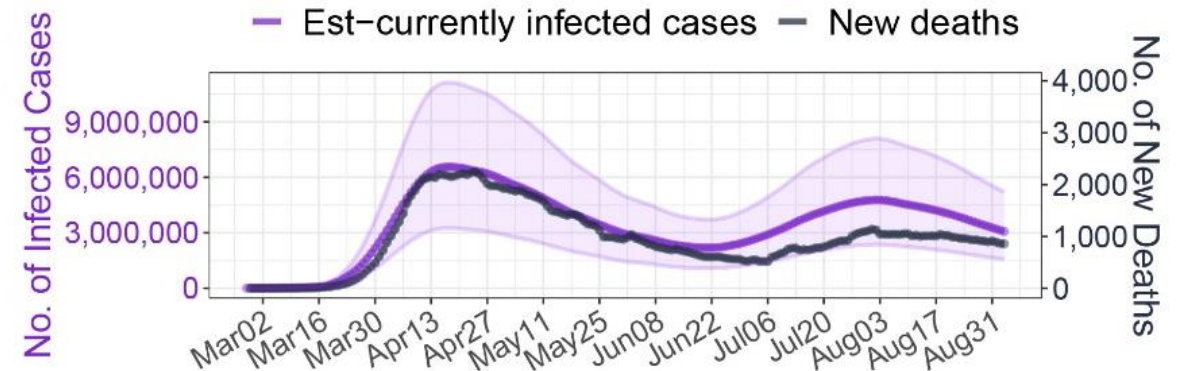
Estimation of the fraction of COVID-19 infected people in U.S. states and countries worldwide

Jungsik Noh ^{*}, Gaudenz Danuser

Estimated Currently Infected Cases

US, as of 2020-09-03:

3,073,341 (0.93% of pop.) [0.47%–1.58%]



Back to dismal 2020...

Mar 31, 2020



The New York Times

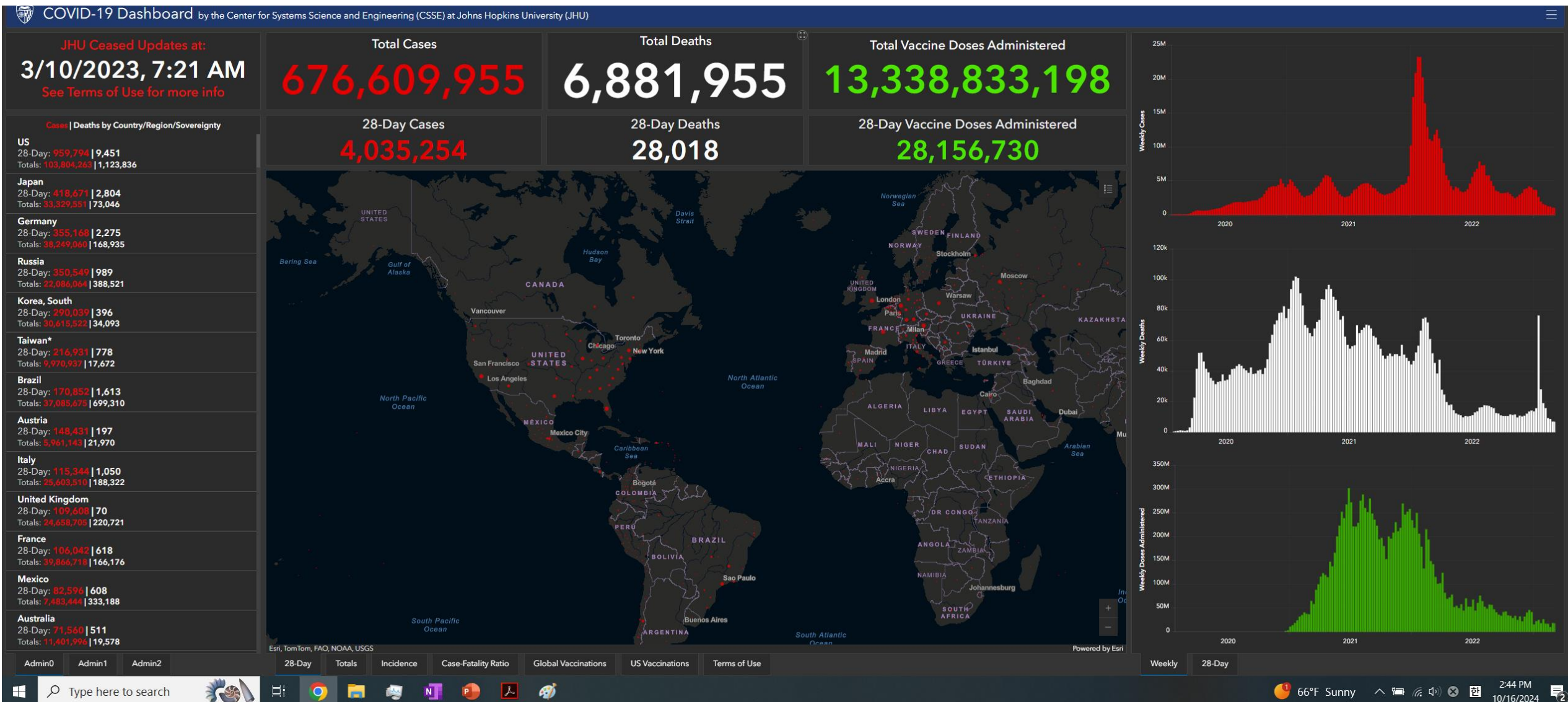


N.Y.C.'s 911 System Is Overwhelmed. 'I'm Terrified,' a Paramedic Says. - The New York Times

Visit >

Data Science for the Pandemic (Data collection)

- Center for Systems Science and Engineering (CSSE) at Johns Hopkins University



Data Science for the Pandemic (Prediction)

May 14, 2020

- INSTITUTE FOR HEALTH METRICS AND EVALUATION (IHME)



Deborah Birx praised for managing coronavirus and White House, Trump

Visit >

Worldwide COVID-19 (12-12-2019 ~ 8-13-2020)

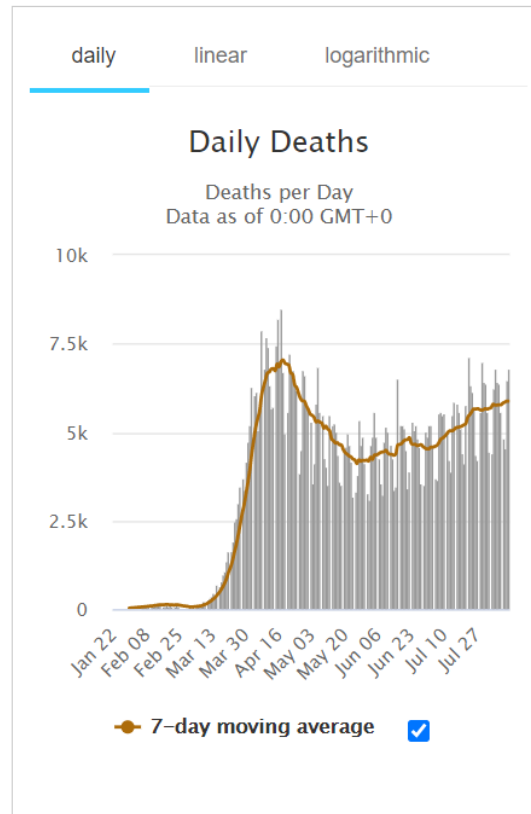
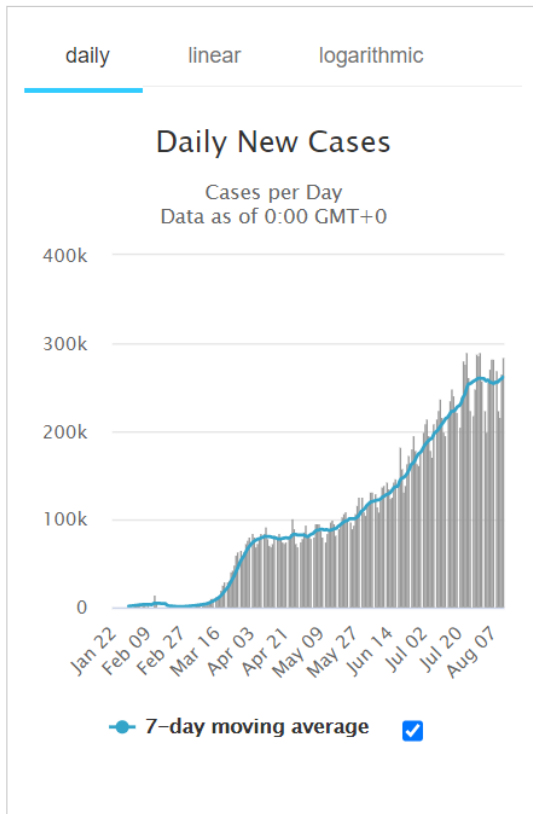
Coronavirus Cases:

20,963,071

[view by country](#)

Deaths:

750,212



What is known
(data)

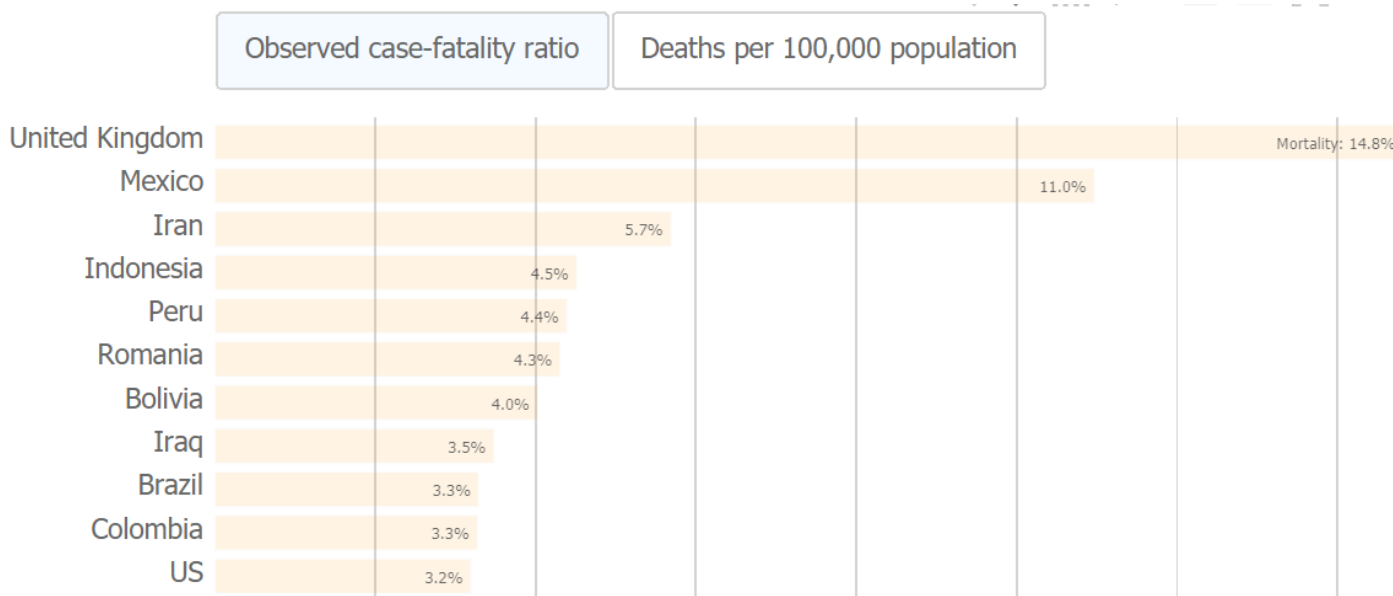
- Daily (reported) new infected cases
- Daily (reported) deaths
- Daily test positivity rates
- Daily hospitalization rates
- ...

What is unknown
(parameters)

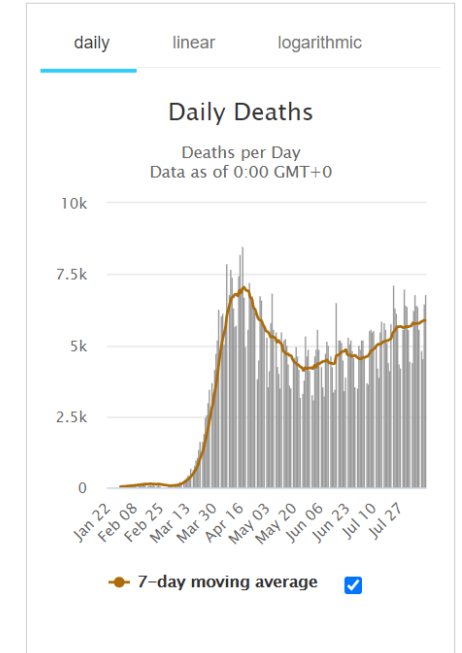
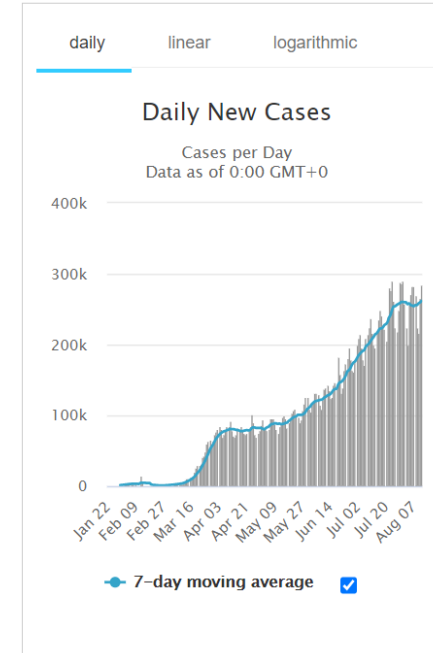
- Key pandemic parameters are unknown
 - Infection-Fatality-Rate (IFR)
 - Time period between:
 - Infection to symptom-onset
 - Infectious period
 - Infection to death
 - Infection to recovery
- ...
- *(Personally) # of currently infected people in Collin county*

Numbers did not make sense: Case-Fatality-Rate (as of Aug 2020)

Spatial variation



Temporal variation



Under-ascertainment due to Low Testing Capacity

Early evidences of substantial undocumented infections

- **Ascertainment rate**

$$= \frac{(\# \text{ of confirmed cases})}{(\# \text{ of actual cases})}$$

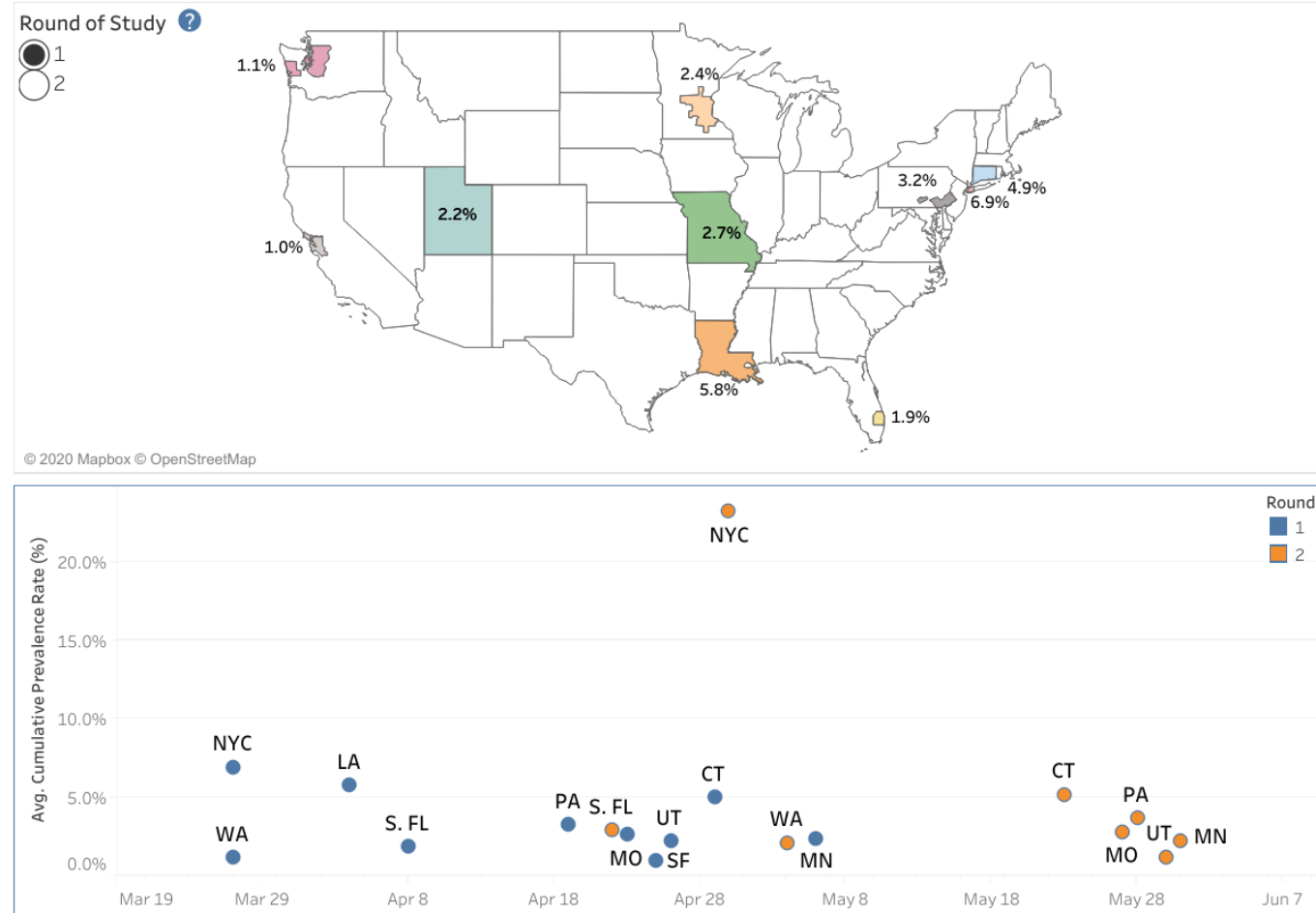
- Using computational modeling, Li et al. (Science, 2020) estimated that only 14% of all infections in China were detected or laboratory-confirmed in January 2020.

- From seroprevalence studies, CDC reported that the ascertainment rates were:
 - 4.2% in MO
 - 8.9% in NY
 - 16.7% in CT until March or April 2020

Seroprevalence Estimates ?

The map shows the seroprevalence estimates for the selected round of study (Round 1 or 2)
The bottom chart shows the seroprevalence estimates for all sites and rounds..

About the study

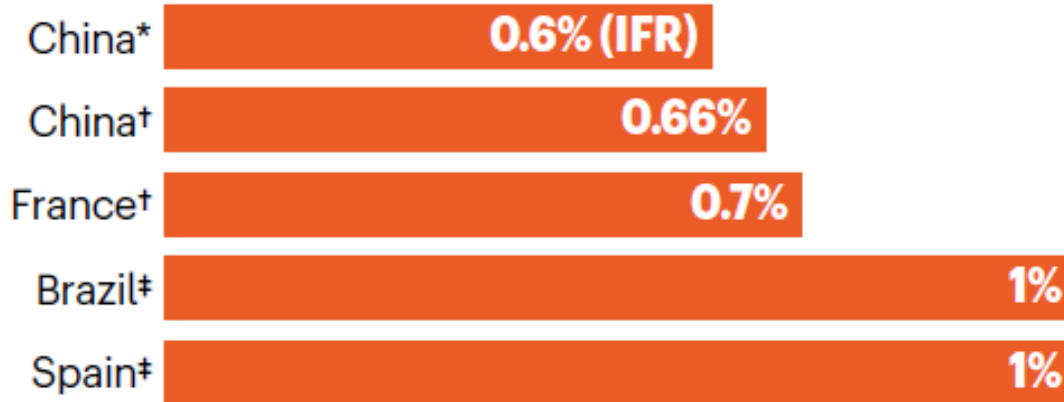


Limitations: A full list of limitations interpreting this data can be found at
<https://www.cdc.gov/coronavirus/2019-ncov/cases-updates/commercial-lab-surveys.html#interpreting-serology-results>

Under-reported Infections and Infection-Fatality-Rate (IFR)

HOW DEADLY IS SARS-COV-2?

The infection fatality rate (IFR) is the proportion of people with COVID-19 who will die from the disease. Estimates are for specific regions, and can vary depending on demographics, health-care access and study methodology.



*Estimate based on natural experiment. †Estimate based on modelling.
‡Estimate based on prevalence data.



United States

Confirmed

5.22M

Recovered

-

Deaths

166K

- Verity et al. (Lancet, March 2020):
 - key pandemic parameters based on early pandemic data

Infection prevalence estimated from international Wuhan residents who were repatriated to their home countries



Age-stratified confirmed cases (>40K) in China during January

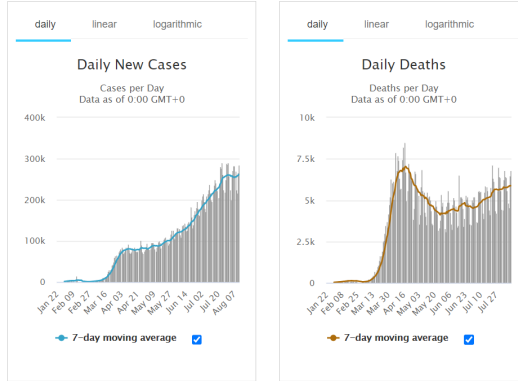


“Our estimated overall infection fatality ratio for China was 0.66% (0.39–1.33)”

- Infection to death: 18 days
- Infection to recovery: 25 days

Formulate a data-science problem

Datasets



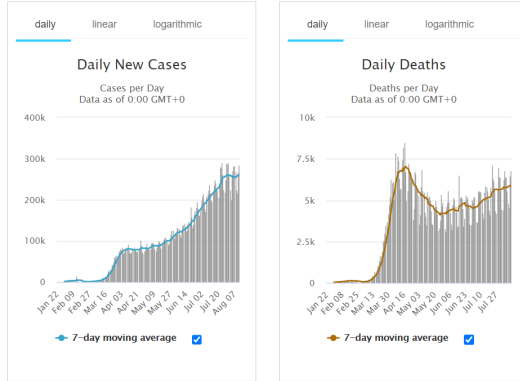
- IFR: 0.66% (0.39–1.33)
- Infection to death: 18 days
- Infection to recovery: 25 days

Goal of a DS project

*# of currently infected
people in Collin county*

Formulate a data-science problem

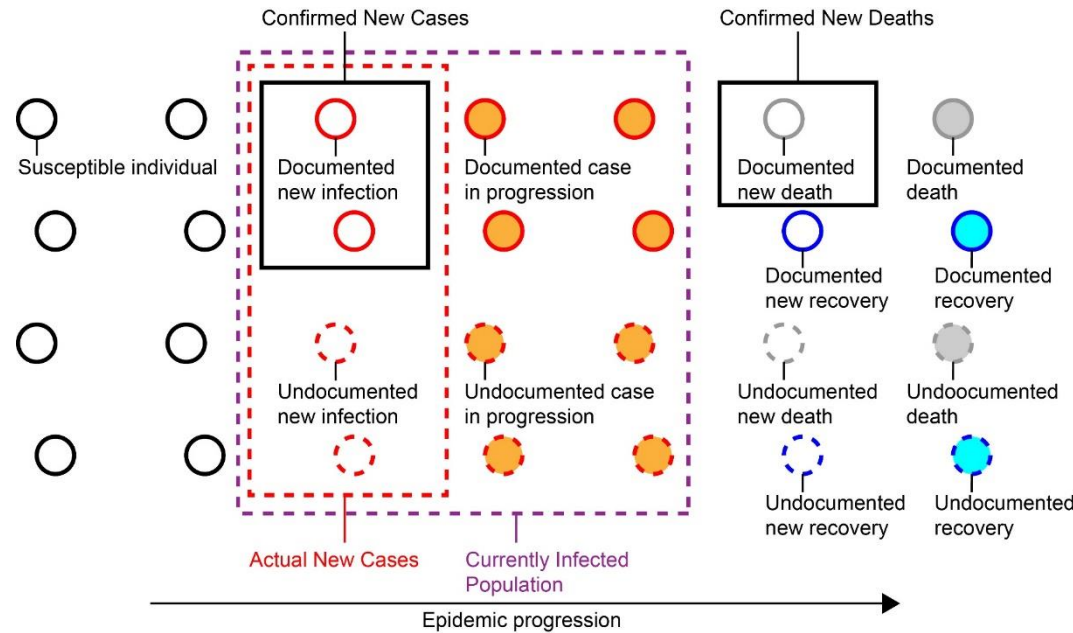
Datasets



- IFR: 0.66% (0.39–1.33)
- Infection to death: 18 days
- Infection to recovery: 25 days

Basic equation in epidemiology

Susceptible(t) \rightarrow Infected(t) \rightarrow Recovered(t) or Death(t)



Goal of a DS project

of currently infected people in Collin county

Total Infected(t)
= Currently Infected(t)
+ Total Death(t)
+ Total Recovered(t)

Assumption

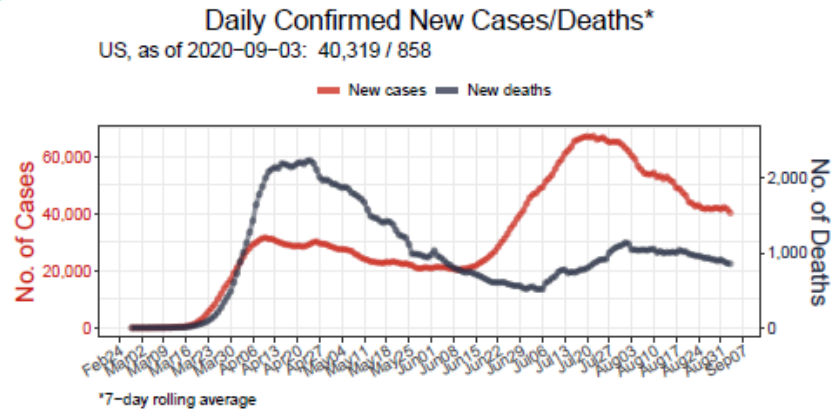
- (A1) Assume that # of deaths is accurate.
(A2) The wide uncertainty of the IFR estimate is expected to cover the true IFRs of many countries and U.S. states.

Daily cases and deaths

1. Actual cases(t - 18) = 1/0.66 * Death(t)
2. Recovered(t+7) = (1/0.66 - 1) * Death(t)
3. Currently Infected(t)
= Total actual cases(t) - Total Death(t) - Total Recovered(t)

Initial latent time series

B



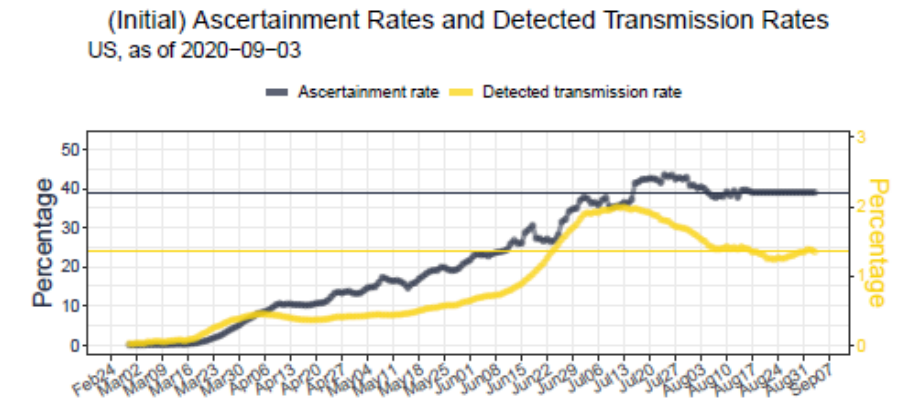
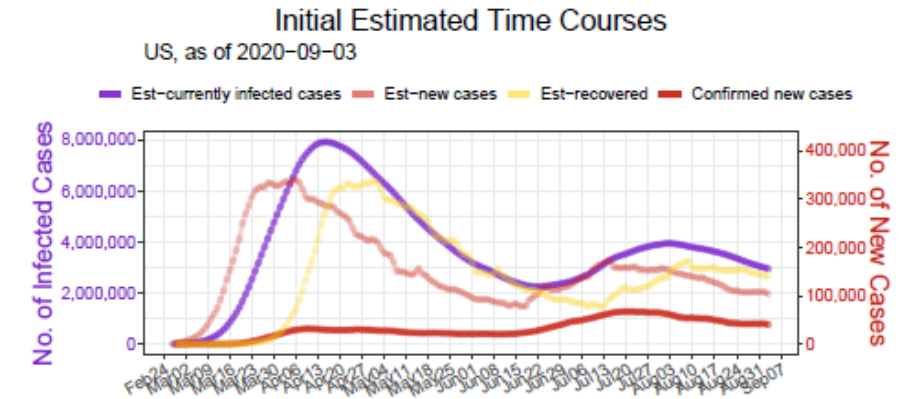
IFR: 0.66% (0.39–1.33)
Infection to death: 18 days
Infection to recovery: 25 days



Initial estimates

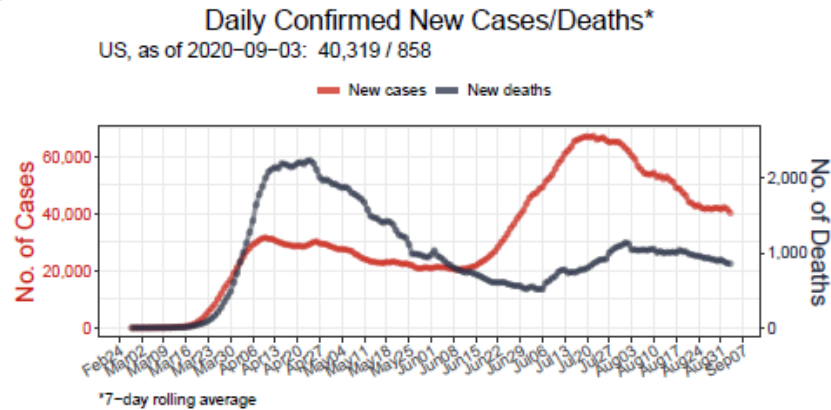
1. **Actual cases**($t - 18$) = $1/0.66 * \text{Death}(t)$
2. **Recovered**($t+7$) = $(1/0.66 - 1) * \text{Death}(t)$
3. **Currently Infected**(t) = Total actual cases(t) – Total Death(t) – Total Recovered(t)
4. Daily ascertainment rate(t)
= Confirmed cases(t) / Actual cases(t)

1. Initialization



Initial latent time series

B



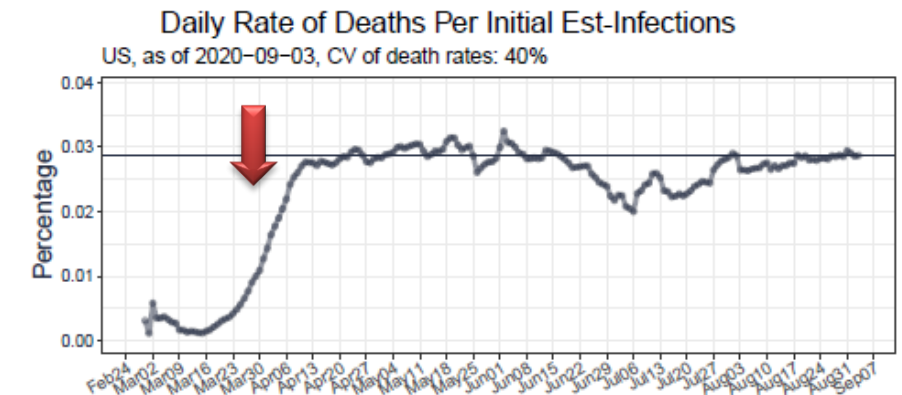
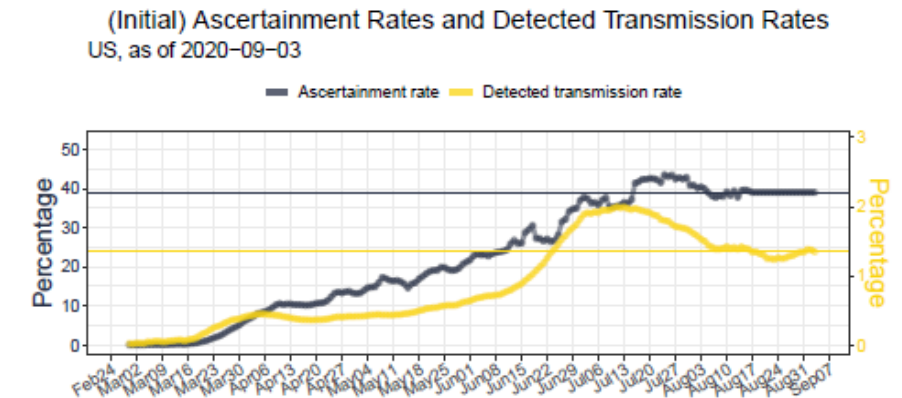
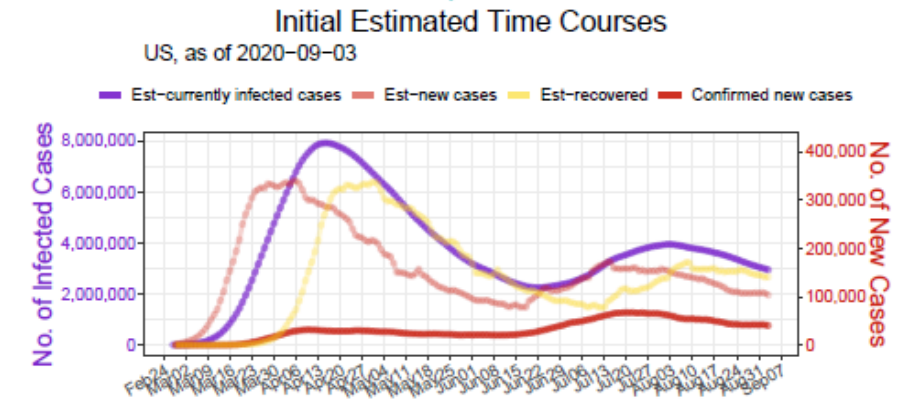
IFR: 0.66% (0.39–1.33)
Infection to death: 18 days
Infection to recovery: 25 days



Initial estimates

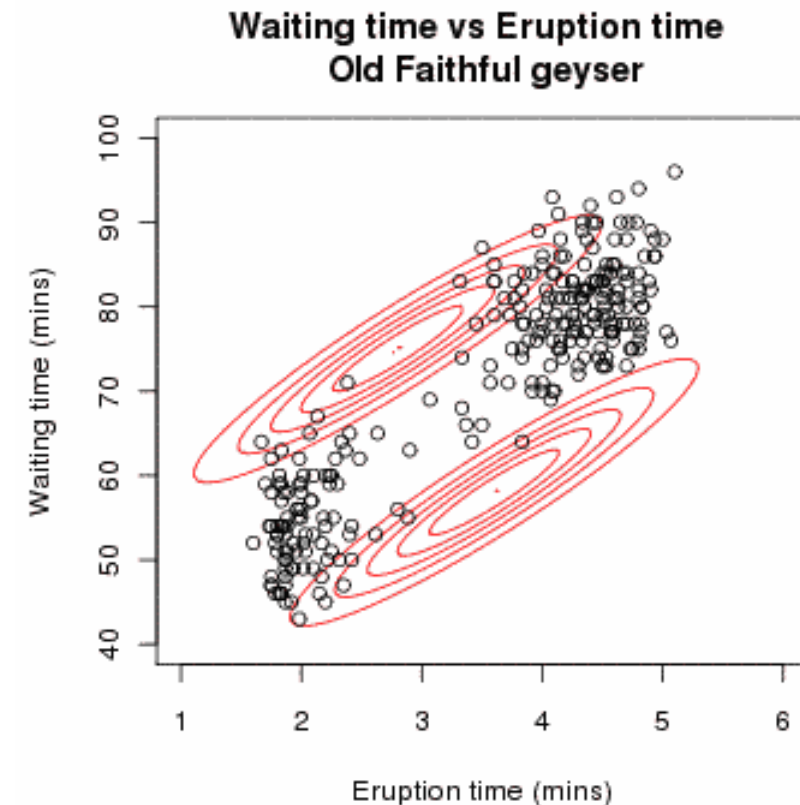
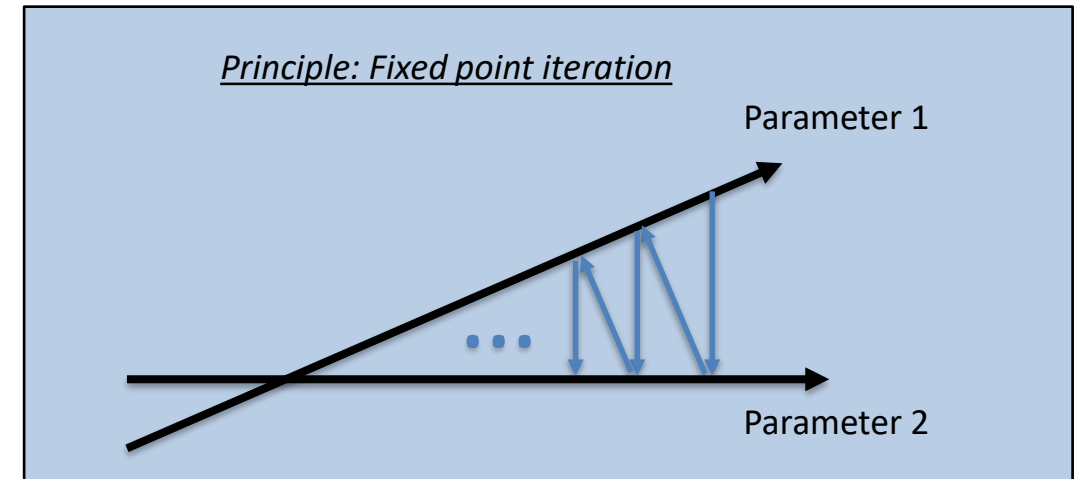
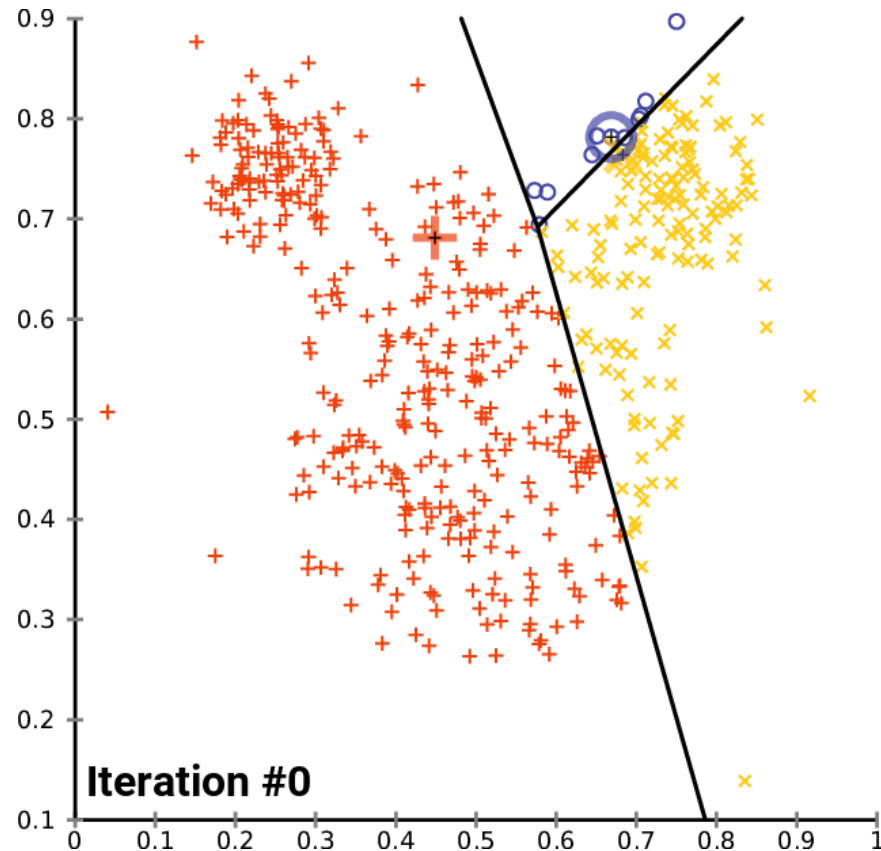
1. **Actual cases**($t - 18$) = $1/0.66 * \text{Death}(t)$
2. **Recovered**($t+7$) = $(1/0.66 - 1) * \text{Death}(t)$
3. **Currently Infected**(t) = Total actual cases(t) – Total Death(t) – Total Recovered(t)
4. Daily ascertainment rate(t)
= Confirmed cases(t) / Actual cases(t)
5. **Detected transmission rate**(t)
= Confirmed cases(t) / Currently Infected($t-1$)
6. Rate of Deaths per Currently Infected(t)
= Deaths(t) / Currently Infected($t-1$)

1. Initialization

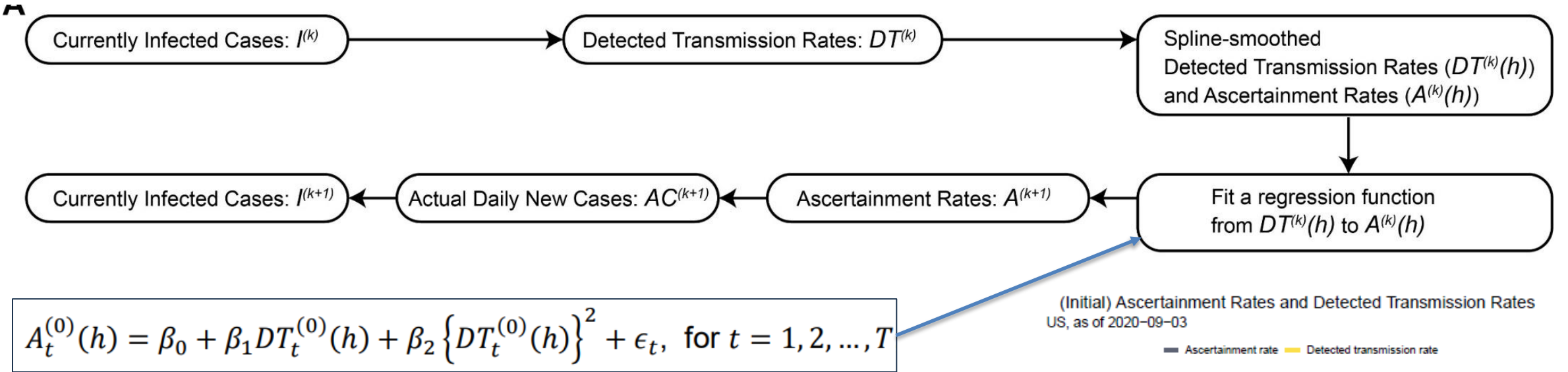


K-means and EM for GMM

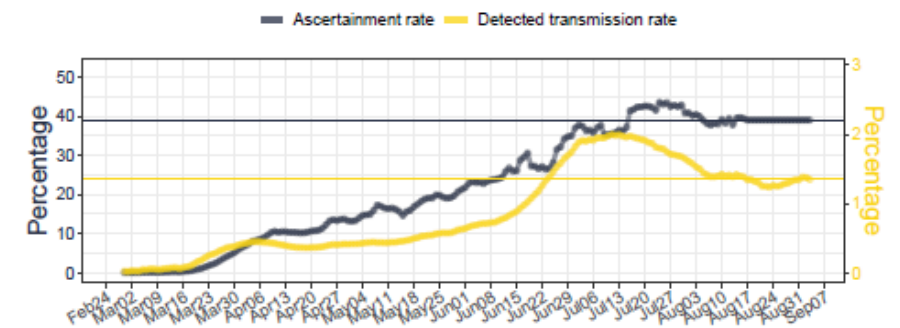
- Both are techniques to find clusters in data points.
- K-means is a hard-thresholding.
- EM for GMM is a soft-thresholding.



EM iterations with COVID time series

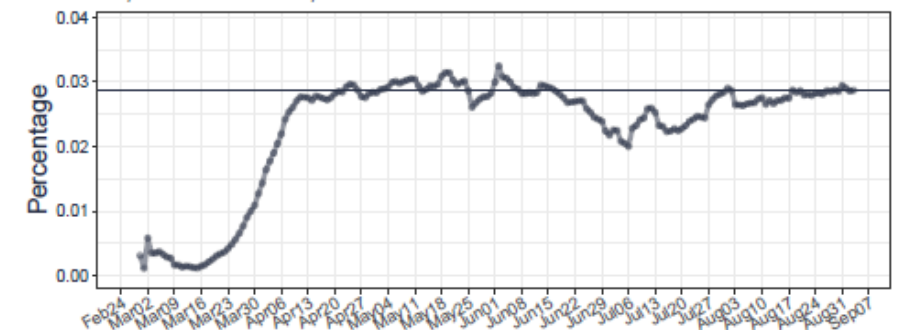


(Initial) Ascertainment Rates and Detected Transmission Rates
US, as of 2020-09-03



Daily Rate of Deaths Per Initial Est-Infections

US, as of 2020-09-03, CV of death rates: 40%

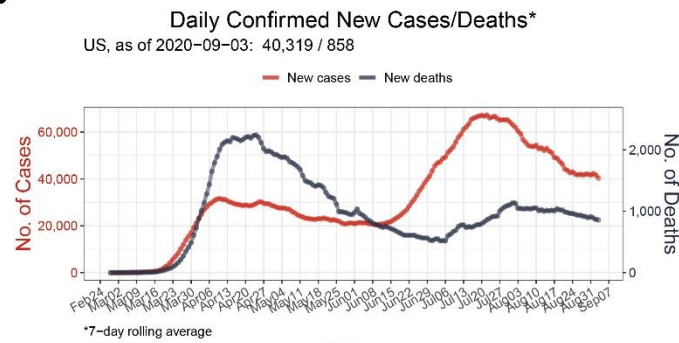


Initial estimates

1. **Actual cases**($t - 18$) = $1/0.66 * \text{Death}(t)$
2. **Recovered**($t+7$) = $(1/0.66 - 1) * \text{Death}(t)$
3. **Currently Infected**(t) = Total actual cases(t) – Total Death(t) – Total Recovered(t)
4. Daily ascertainment rate(t)
 $A_t = \text{Confirmed cases}(t) / \text{Actual cases}(t)$
5. Detected transmission rate(t)
 $DT_t = \text{Confirmed cases}(t) / \text{Currently Infected}(t-1)$
6. Rate of Deaths per Currently Infected(t)
= Deaths(t) / Currently Infected($t-1$)

Converged estimates

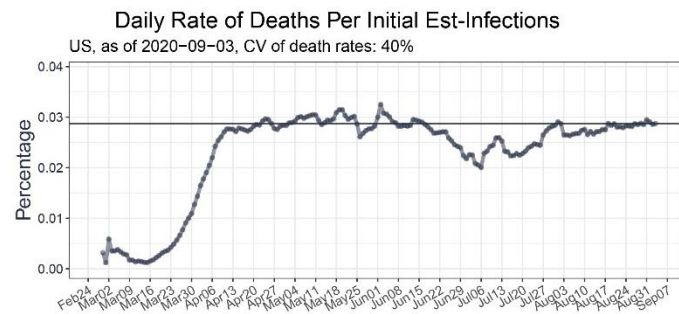
B



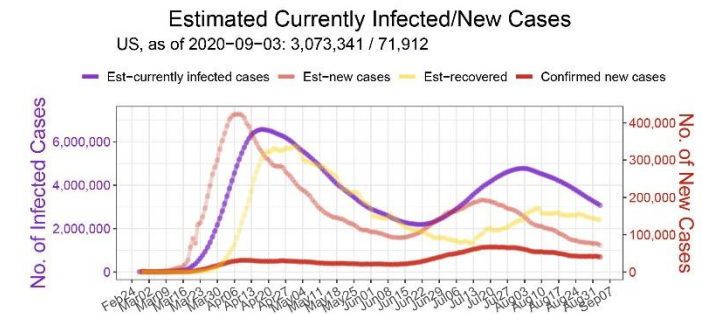
1. Initialization



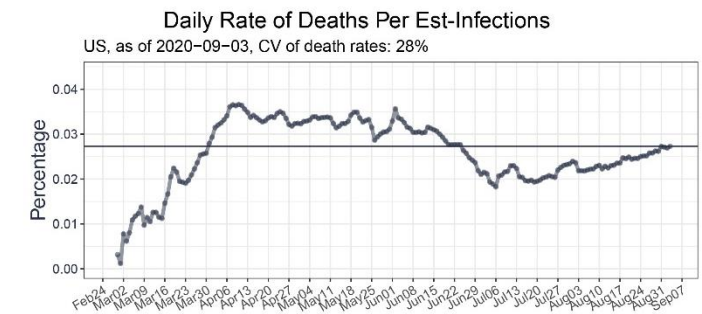
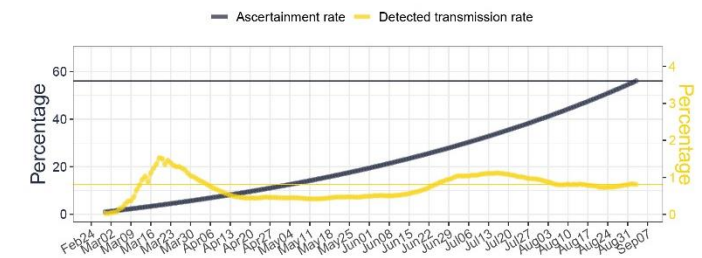
(Initial) Ascertainment Rates and Detected Transmission Rates
US, as of 2020-09-03



2. EM iterations



(Converged) Ascertainment Rates and Detected Transmission Rates
US, as of 2020-09-03, splinePar: 1.6

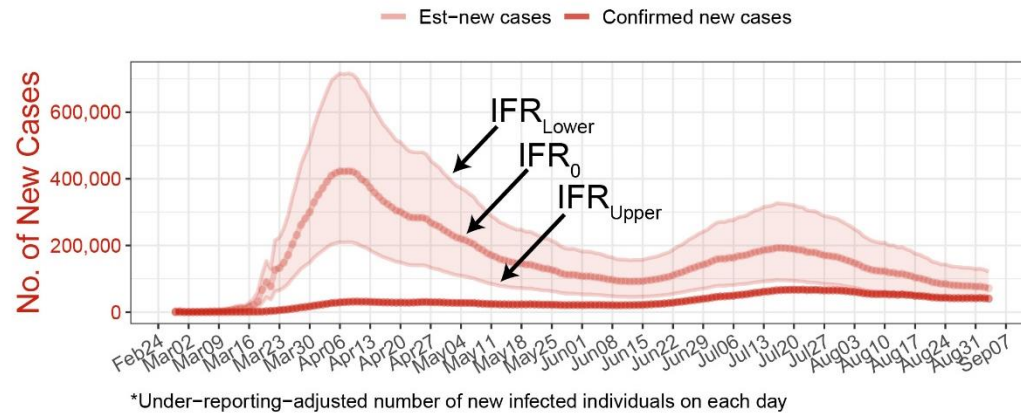


Confidence intervals

3. Calculate Confidence Intervals

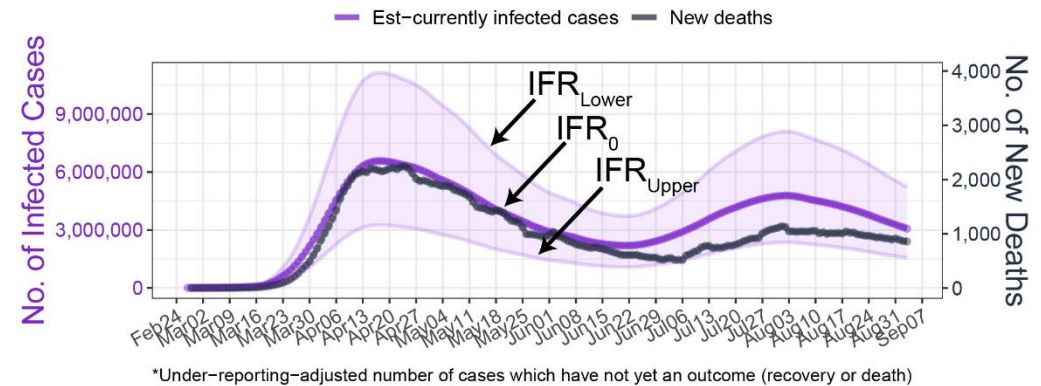
Estimated*/Confirmed New Cases

US, as of 2020-09-03: 71,912 / 40,319 [40,726–121,698]



Estimated Currently Infected Cases*

US, as of 2020-09-03:
3,073,341 (0.93% of pop.) [1,561,253–5,207,578] [0.47%–1.58%]

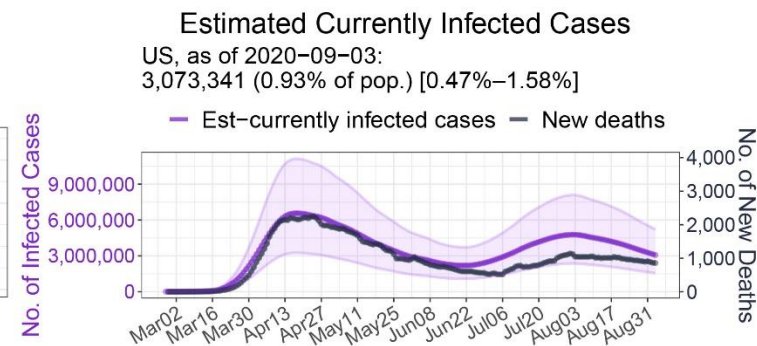
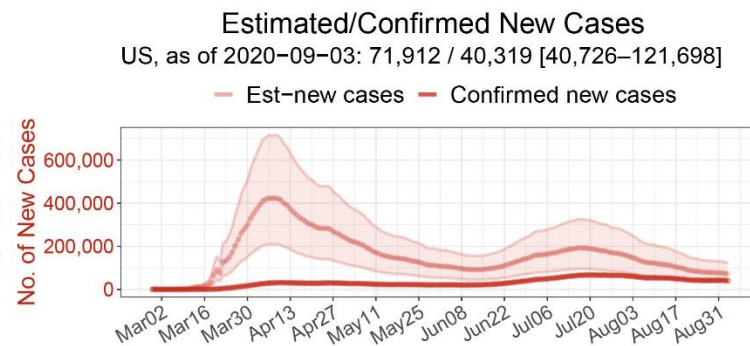
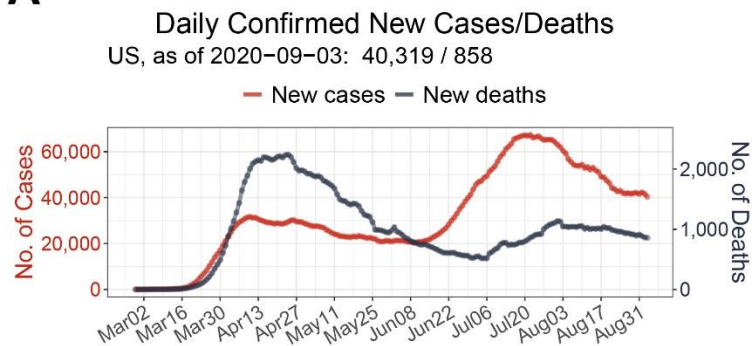


- IFR: 0.66% (0.39–1.33)
- Infection to death: 18 days
- Infection to recovery: 25 days

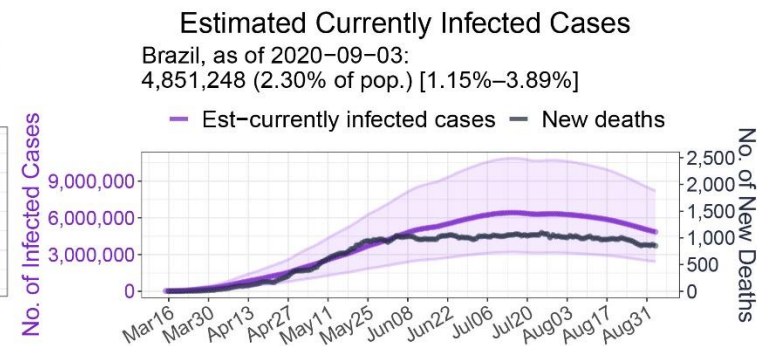
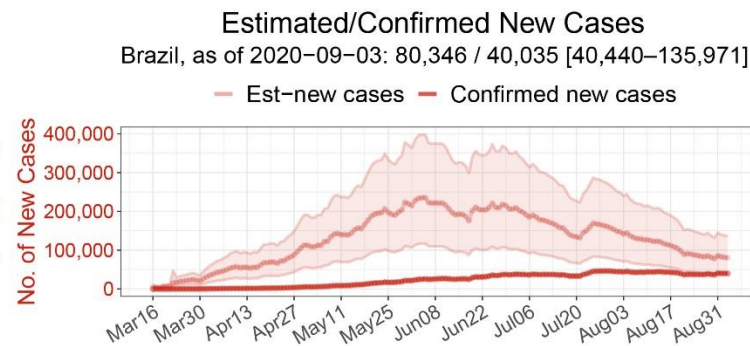
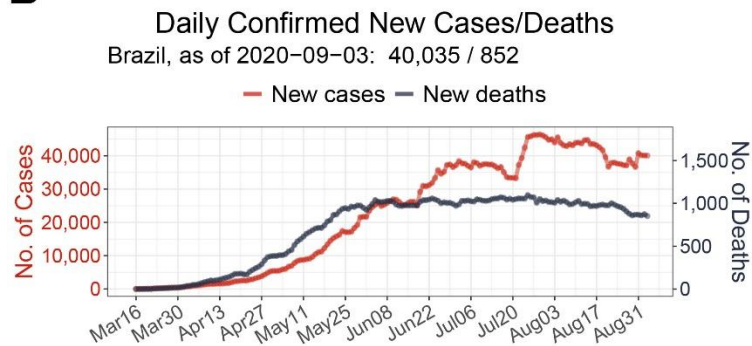
Actual cases in US, Brazil and Louisiana

Goal of a DS project

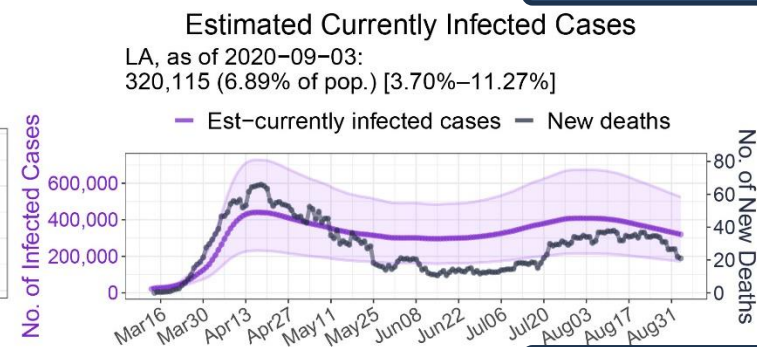
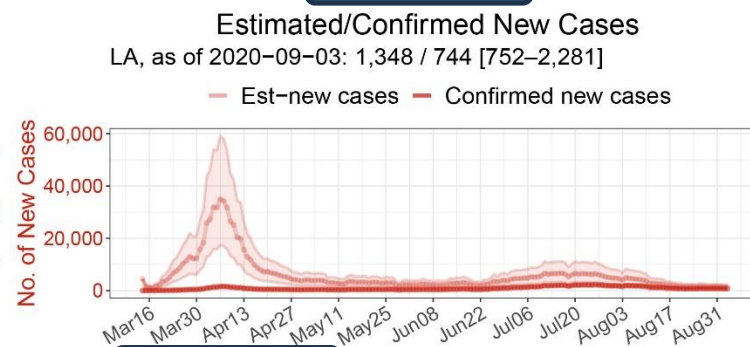
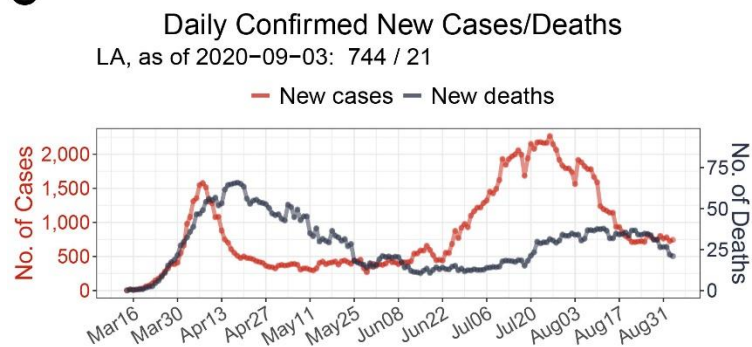
A



B



C



$A(t) < 10\%$

$I(t)$ rate $\approx 1\%$

$A(t) \approx 10\%$

$I(t)$ rate $\approx 2.3\%$

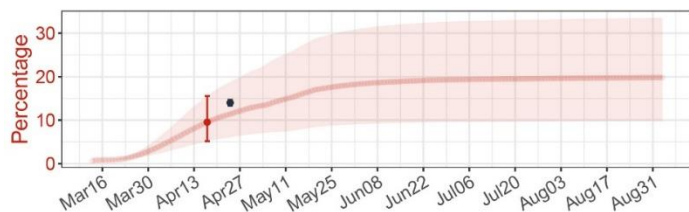
$A(t) \approx 5\%$

$I(t)$ rate $\approx 7\%$

Validation using seroprevalence survey data

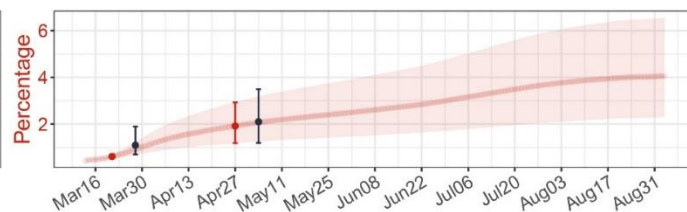
NY, as of 2020-09-03: 19.8% [9.8%–33.5%]

• Est-cumulative incidence • Seroprevalence



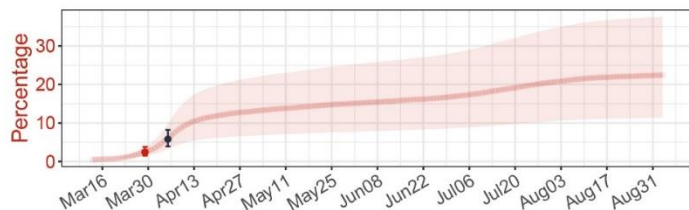
WA, as of 2020-09-03: 4.1% [2.3%–6.5%]

• Est-cumulative incidence • Seroprevalence



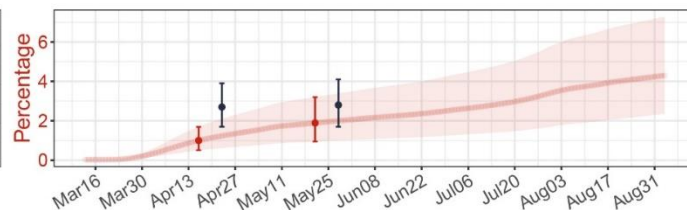
LA, as of 2020-09-03: 22.4% [11.4%–37.6%]

• Est-cumulative incidence • Seroprevalence

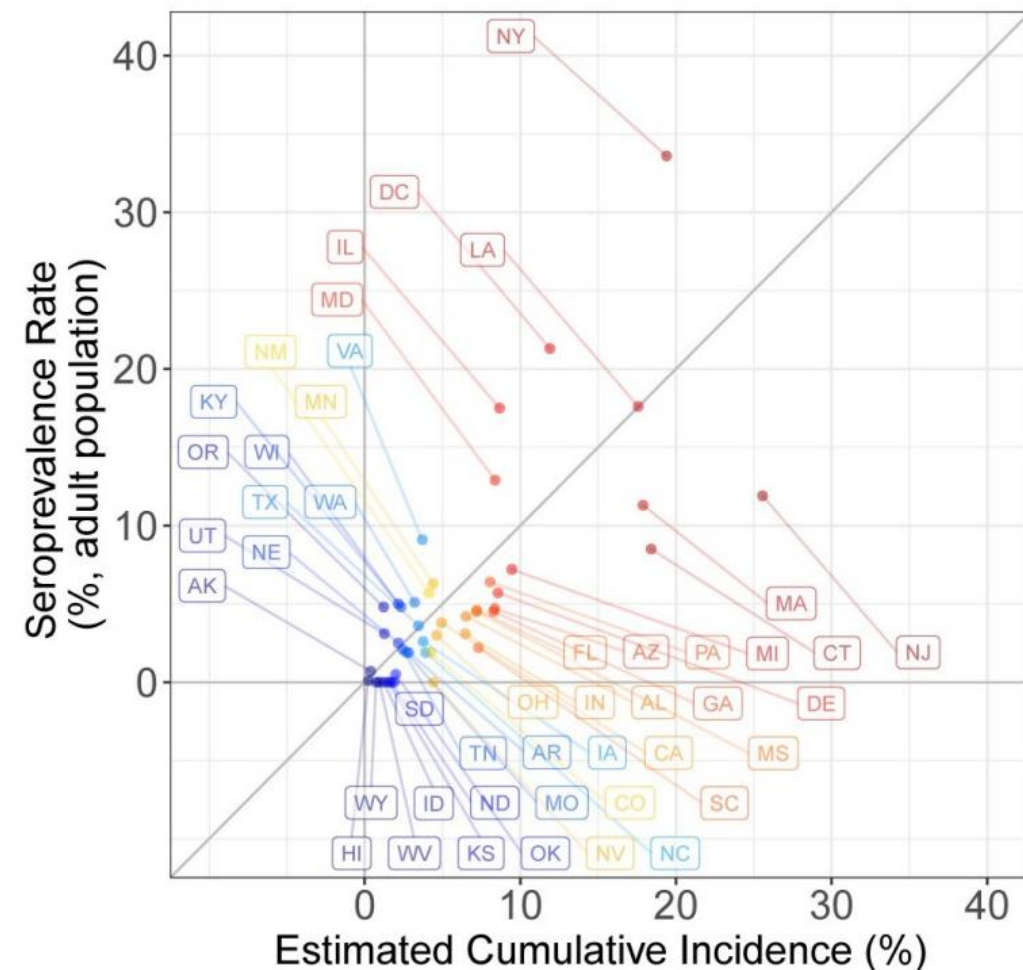


MO, as of 2020-09-03: 4.3% [2.4%–7.3%]

• Est-cumulative incidence • Seroprevalence



CDC survey



Nationwide blood tests with
dialysis patients (n = 28,503)
during July 2020

Summary

1. Goal of a DS project

2.

