# Modelling COVID-19

## Making the best of bad data

---

Luiz Max Carvalho [lmax.fgv@gmail.com]

**FGV EMAp**

# Acknowledgements

1

## Data

Why is the data "bad"?

## Data

Why is the data "bad"?

## Models

Are all models wrong? Are some useful?

## Data
Why is the data "bad"?

## Models
Are all models wrong? Are some useful?

## Some results
What can semi-mechanistic and non-mechanistic models do for us?

## Data

Why is the data "bad"?

## Models

Are all models wrong? Are some useful?

## Some results

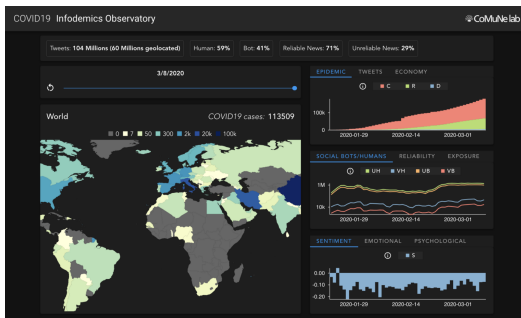What can semi-mechanistic and non-mechanistic models do for us?

## Future

What have we learned? How can we improve going forward?

◎ Lots of wonderful efforts to collect and make sense of data:

- Johns Hopkins University;
- CoronaNet data collection project;
- Nextstrain;
- Covid19 Infodemics Observatory;
- Wesley Cota's website;
- Brasil.IO;
- Infogripe.

- ◎ Cases per day;
- ◎ Deaths per day;
- ◎ Interventions: what, when, for how long;
- ◎ Tests per day;
- ◎ Mobility: Google, Apple, In Loco reports;
- ◎ Serological surveys;
- ◎ Genomic information on the virus.

◎ Inconsistent criteria;

## Part 1: Data problems

- ◎ Inconsistent criteria;
- ◎ Reporting (notification) delays;

## Part 1: Data problems

◎ Inconsistent criteria;

◎ Reporting (notification) delays;

◎ Underreporting;

◎ Inconsistent criteria;
◎ Reporting (notification) delays;
◎ Underreporting;
◎ Completeness;
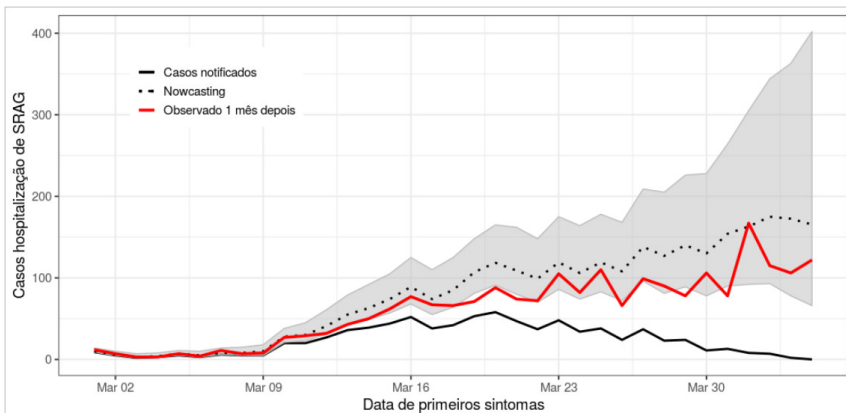
◎ Case definition may vary;
  Here is SRAG:

  - Fever (above 37,8 Celsius) <u>AND</u>;
  - Cough OR sore throat <u>AND</u>;
  - Difficulty breathing OR Dyspnea OR $O_2$ saturation below 95% <u>AND</u>;
  - Was hospitalised OR died with these symptoms.

◎ Testing is inconsistent and uncertain;

  - RT-PCR and antibody-detecting tests vary in scope and applicability;
  - Sensitivity and specificity.
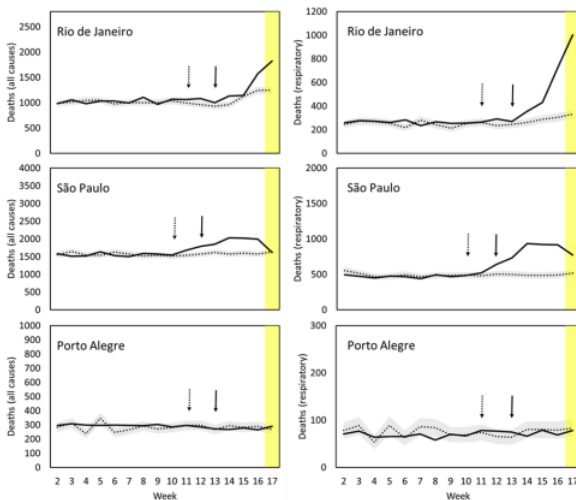
# Data problems II: Delay

**Important:** delayed data <u>exists</u>, it just has not come into the system yet. See Bastos et al. (2019).

Cases that have not come into the system and probably won't.

◎ One approach is excess mortality (e.g. Freitas et al., (2020))

Case data is usually incomplete/inconsistent.

◎ The beauty of open fields...

First, an incomplete model taxonomy (Ronald Ross [1857-1932]):

- ◎ *A priori* pathometry;
  - ○ Representation of (disease) mechanisms;
  - ○ Systems of differential equations;
  - ○ "Mathematical models".
- ◎ *A posteriori* pathometry;
  - ○ Curve fitting;
  - ○ "Statistical models".

# A (fancy) SEIR model

$$\frac{dS}{dt} = -\lambda[(1 - \chi)S],$$

$$\frac{dE}{dt} = \lambda[(1 - \chi)S] - \alpha E,$$

$$\frac{dI}{dt} = (1 - p)\alpha E - \delta I - \phi I,$$

$$\frac{dA}{dt} = p\alpha E - \gamma A,$$

$$\frac{dH}{dt} = \phi I - (\rho + \mu)H,$$

$$\frac{dR}{dt} = \delta I + \rho H + \gamma A,$$

$$\lambda := \beta(I + A).$$



SEQIAHR model

◎ Hard to fit to data;

◎ Usually quite sensitive to "minor" features of the data;

◎ Projections can be thrown off as a result;

◎ Hard to integrate various sources of information.

## Statistical models

- More flexible (splines, GAMs, Gaussian processes);
- More natural accommodation of stochasticity and uncertainty;
- Less insight into mechanisms;
- More difficult to simulate scenarios (but not impossible!).

## General considerations on modelling

- ◎ Meme: "All models are wrong, some are useful" (Box, 1976);
- ◎ "Many quotes are cool, few are universal";
- ◎ Usefulness is context-dependent $\implies$ goal-oriented modelling;
- ◎ Every modelling endeavour will sacrifice "realism" for tractability, the important consideration is what to leave out *relative* to the task at hand (prediction, estimation, scenario modelling).

Many projects devoted to predicting the numbers of cases and deaths over time.

◎ UFMG;

◎ PUC-Rio;

◎ UFV;

◎ BRAM-COD (predicts ICU occupation as well);

For the first example, we will be discussing the of Kubinec & Carvalho (2020).

◎ If we have covariate data per state, say, can we study their effects on epidemic progression?;

◎ Can we make useful predictions even when the model lacks explicit dynamics?

## Setup and notation

Here we will consider an *empirical*, *retrospective* model for
infection rates. Consider time points $t = 1, \ldots, T$ and regions
$c \in C$. We will deal with

$$f_t \left( \frac{I(t)}{S(t) + R(t)} \right)$$

where $f_t : (0, \infty) \to (0, \infty)$ is a historical time trend, which we
will call the "empirical" trend. This model can be seen as local
linear approximation of the $I_c(t)$ curve and cannot be used for
future predictions, for example.

Additionally, we do not observe $I_{ct} := I_c(t)$ directly, but rather
the numbers of tests, $q_{ct}$, and cases, $a_{ct}$, along with a set of
covariates $X_{ct}$.

## Model details

$$I_{ct} \sim \text{Beta}(\alpha_1 + \beta_{O1} \sum_{c=1}^{C} \mathbf{1}(a_{ct'} > 0) + \beta_{S1} \mathbf{X}_{ct} +$$

$$\beta_{I1} t_o + \beta_{I2} t_o^2 + \beta_{I3} t_o^3, \phi)$$

$$q_{ct} \mid \text{pop}_c \sim \text{Beta-Binomial}(\text{pop}_c, \text{logit}^{-1}(\alpha_2 + \beta_q I_{ct}), \phi_q),$$

$$a_{ct} \mid q_{ct} \sim \text{Beta-Binomial}(q_{ct}, \text{logit}^{-1}(\alpha_3 + \beta_a I_{ct}), \phi_a).$$

$$\beta_a \sim \text{Exponential}(.1),$$
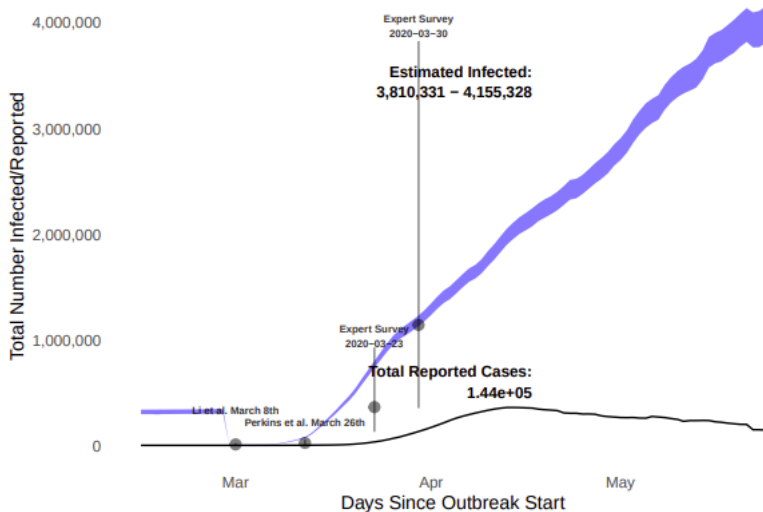$$\beta_{qc} \sim \text{Exponential}(\sigma_q),$$
$$\sigma_q \sim \text{Exponential}(.1),$$
$$\beta_{Si} \sim \text{Normal}(0, 2),$$
$$\beta_{Ii} \sim \text{Normal}(0, 5),$$
$$\alpha_i \sim \text{Normal}(0, 10), i = 1, 2, 3$$

Example I results: reconstructions of $I_t$

# Example I results: covariate effects



Cumulative Effect on Proportion Infected

| | |
|---|---|
| Trump Vote Share | 0.004222% |
| Cardiovascular Deaths | 0.003106% |
| PM 2.5 | 0.000338% |
| Public Health Funding | 0.000245% |
| GDP | −0.000058% |
| % Foreign−Born | −0.002721% |
| No. Providers | −0.003825% |
| % Population <18 | −0.005417% |
| % Smokers | −0.007586% |

Effect on Proportion Infected

# Example I results: testing rates across US states



Additional Proportion Tested of Population for Every Percent Increase in Infected

22

## Limitations

◎ Confounding structures:



◎ "Identifying" prior assumes constant proportion of (under)detection;

◎ Assumes same temporal trend for all locations, which might be unrealistic.

## Example II: where did the virus come from?

In this example, we will look into the findings of Candido et al. (2020).

- ◎ Effective reproductive number ($R_t$) of COVID-19 in major cities;
- ◎ Viral dispersal across the country;

# Example II: $R_t$

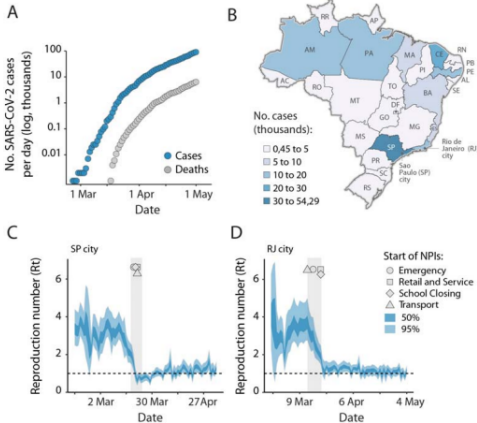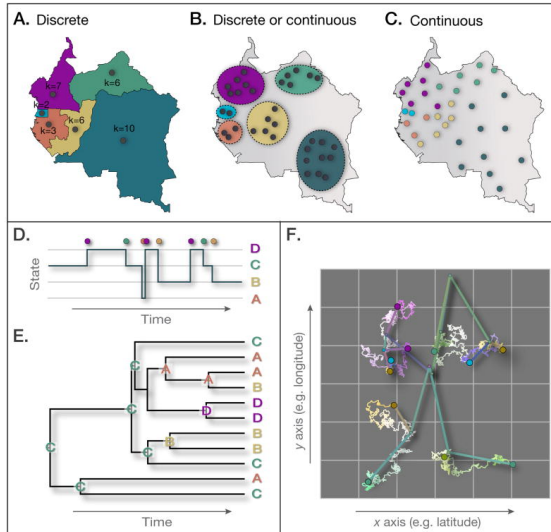$$R_{tm} = R_{0m}\left(2\,\text{logit}^{-1}\left(-\sum_{k=1}^{4}(\alpha_k + \beta_{mk})I_{ktm} + B_k\right)\right)$$

$$P(t) = \exp(Qt)$$ (Faria et al., 2011)

A. First epidemic phase / Second epidemic phase

Population density per sq. Km:
0 – 2
2 – 11
11 – 24
24 – 48
48 – 216
216 – 14804

Timescale of internal nodes:
15 Feb – 18 Mar
18 Mar

18 Mar – 15 Apr

Dispersal direction:

B. No. observed migratoin events — within state / among state — start of NPIs — $T_{shift}$ — Date (months and days of 2020): 15 Feb, 15 Mar, 15 Apr

C. Average dailiy distance travelled per passenger (km, log) — start of NPIs — Date (months and days of 2020): 15 Feb, 15 Mar, 15 Apr

## Part 4: What have we learned?

After this bird's-eye tour, what have we gathered?

◎ The first truly "data-driven" pandemic has exposed the inadequacy of our reporting protocols;

◎ On the other hand, it has also shown the incredible potential for rapid sharing of data and methods and for global collaboration;

◎ Whilst we would like to understand *mechanisms*, semi-structured statistical models can help answer questions about

  ○ Factors associated with infection rates in each location;
  ○ Indicators of transmission risk such as $R_t$;
  ○ Disease dispersal patterns.

## Practical recommendations

This is our dry-run for the end of the World. A much bigger and deadlier epidemic **will** come. How we prepare will tell whether we thrive or whither away as a species.

- ◎ Government: strengthen access to health care (Werneck & Carvalho, 2020);
- ◎ Government/funding agencies: fund surveillance programmes;
- ◎ Funding agencies: fund risk management projects targeting health crises;
- ◎ Modellers: research adequate methods for accommodating uncertainty;
- ◎ Modellers: research model-ensemble methods;

## Data issues

Inconsistent criteria, delays, underreporting, incompleteness

### Data issues

Inconsistent criteria, delays, underreporting, incompleteness

### "Empirical" models can help answer pressing questions

Which factors are associated with infection rates? Where did the virus come from and where did it go?

### Data issues

Inconsistent criteria, delays, underreporting, incompleteness

### "Empirical" models can help answer pressing questions

Which factors are associated with infection rates? Where did the virus come from and where did it go?

### Fund surveillance programmes

Accurate and timely monitoring are a must in ever more connected world.

# THE END