

Statistical methods for epidemiological surveillance

Leo Bastos (PROCC/Fiocruz) and Luiz Max Carvalho (FGV/EMAp)

July 5, 2022

Outline

- ① Introduction to infectious disease epidemiology
- ② Mathematical/deterministic models
 - SIR-like models
 - Semi-structured models
- ③ Real time analyses
 - Nowcasting
 - Forecasting infectious diseases

- 1 Introduction to infectious disease epidemiology
- 2 Mathematical/deterministic models
- 3 Real time analyses

(Infectious) Disease process

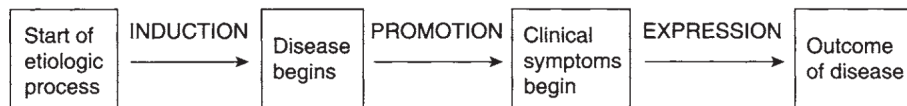


Figure 1.1 *Schematic of disease evolution.*

- ① Start of etiologic process
 - Ex: Infection through a mosquito bite (dengue fever)
- ② Disease begins
 - Viral replication
- ③ Clinical symptoms
 - Fever, headaches,...
- ④ Outcome of disease
 - Cure, hospitalization, death

Epidemiological study of Disease process

- ① Start of etiologic process
 - Can we avoid infection?
- ② Disease begins
 - Is our immune system prepared?
- ③ Clinical symptoms
 - Can we treat or avoid evolution?
- ④ Outcome of disease
 - How to reduce the burden?

Infectious disease epidemiology (IDE)

- Infectious disease epidemiology (IDE) is the study of how and why infectious diseases emerge and spread among different **populations**, and what strategies can prevent or contain the spread of disease at the population level.
- Why is this important?



Data type: In Infectious Disease Epidemiology

- Infectious disease data is mainly binary

$$Z_i = \begin{cases} 1 & \text{person } i \text{ is infected with pathogen A or has a disease D,} \\ 0 & \text{otherwise.} \end{cases}$$

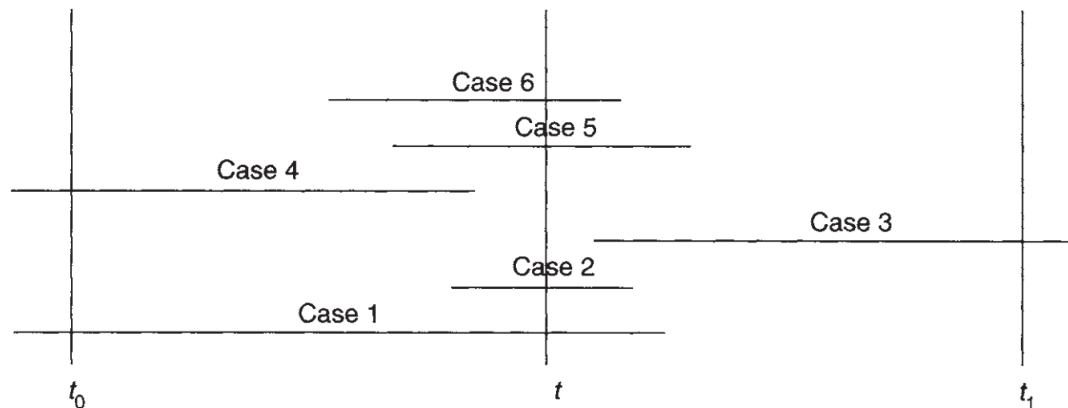
- Counts of cases through time

$$Y_t = \sum_i Z_{i,t}$$

- Induce time and spatial dependence is important, since dependence is present by definition of infectious diseases

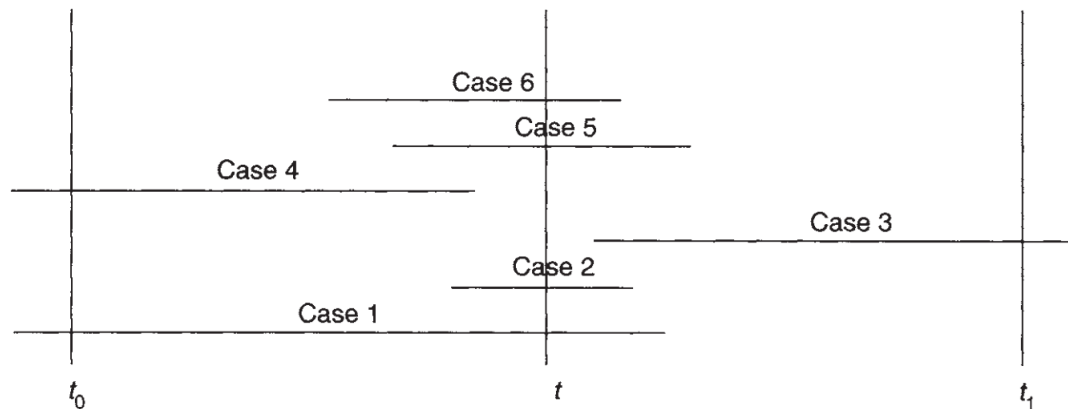
Important measures of disease occurrence

- Disease prevalence and incidence both represent proportions of a population determined to be diseased at certain times.
- Suppose there are 100 people being followed.



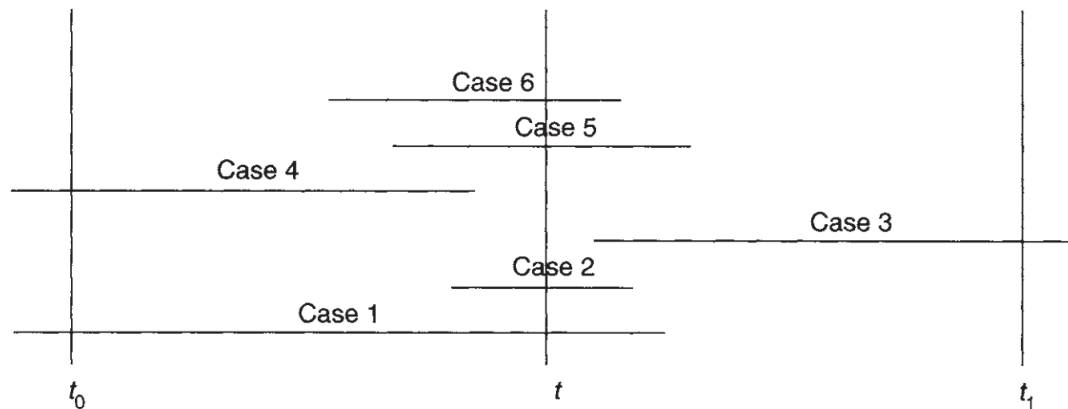
Important measures of disease occurrence

- **Disease prevalence** at time t : $4 / 100$ or $4 / 99$

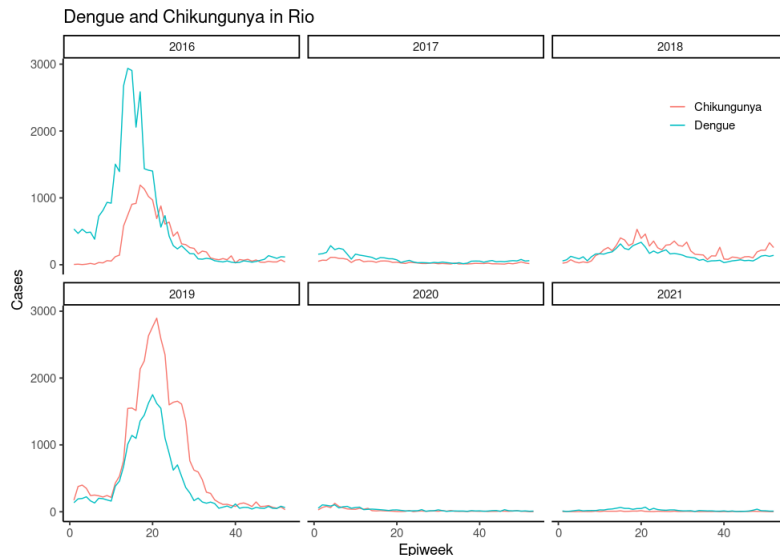


Important measures of disease occurrence

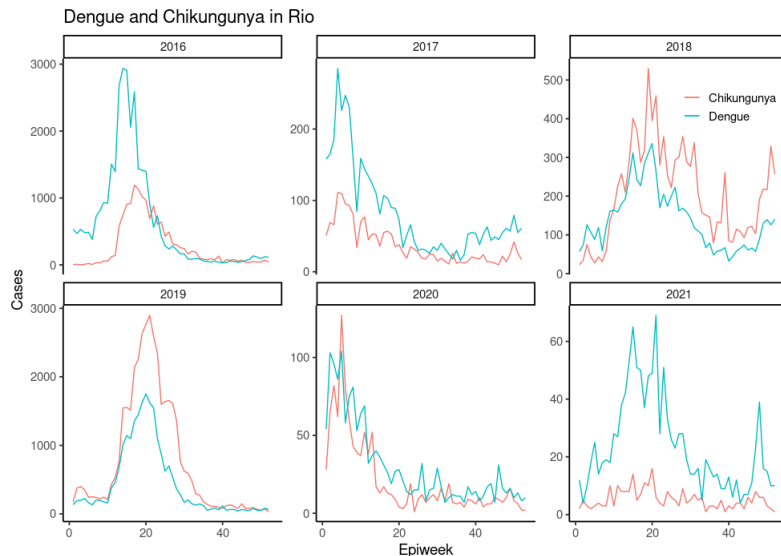
- **Disease incidence** (new cases) in time $[t_0, t_1]$: 4 / 98



Dengue and Chikungunya new cases in Rio de Janeiro



Dengue and Chikungunya new cases in Rio de Janeiro



Outbreak, epidemic, pandemic, endemic

- **Outbreak:** an unexpected increase in the number of disease cases in a small area.
 - Ex: high number of cases of ILI at Vidigal, Rio de Janeiro (Nov/21).
- **Epidemic:** an unexpected increase in the number of disease cases in a specific geographical area.
 - Ex: high number of cases of ILI in several neighbourhoods of Rio de Janeiro (Dec/21).
- **Pandemic:** Exponential growth of the number of cases; Several outbreaks/epidemics spread around the globe.
(Influenza/H1N1, COVID-19)
- **Endemic:** When a disease is consistently present but limited to a particular region..
(Malaria in Brazil's North region)

Modelling outbreaks

- Counting data time series models are important tools to model outbreaks and epidemics.
- the usual models have a low predictive power unless there is some knowledge about the disease dynamics, and if there is historical data (a problem for emerging diseases)
- Endemic diseases with a large data history are those possible to forecast. The long term forecasts could be seen as the expected behaviour (real time observed values greater than the expected values is a strong suggestion of an outbreak).

Original Article

Influenza surveillance in Europe: establishing epidemic thresholds by the Moving Epidemic Method

^aDirección General de Salud Pública, Consejería de Sanidad, Valladolid, Spain. ^bThe Radboud University Nijmegen Medical Centre, Nijmegen, The Netherlands. ^cEuropean Centre for Disease Prevention and Control, Stockholm, Sweden. ^dDivision of Health Security, Infectious Diseases and the Environment, WHO Regional Office for Europe, Copenhagen, Denmark. ^eCentro Nacional de Gripe de Valladolid, Universidad de Valladolid, Spain. ^fInstituto Nacional de Saude Doctor Ricardo Jorge, Lisboa, Portugal.

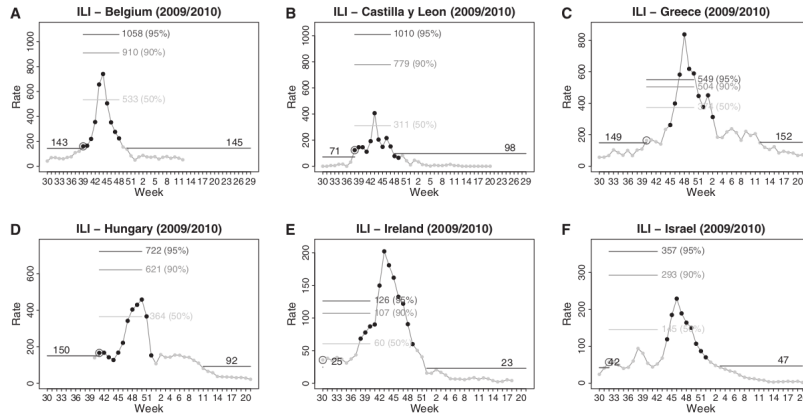
Correspondence: Tomás Vega, Dirección General de Salud Pública, Consejería de Sanidad, Paseo de Zorrilla, 1. 47071 Valladolid, Spain. E-mail: vegaloto@jcyl.es

Accepted 27 June 2012. Published Online 16 August 2012.

MEM: Moving Epidemic Method

- Developed by Vega et al. (2012) focused on Influenza
- MEM is a three-step algorithm:
 - ① Determine start, duration, and end of the epidemic period;
 - ② Estimate epidemic thresholds and epidemic channels defining pre- and post-epidemic levels.

Epidemic thresholds



Epidemic thresholds

R: The Moving Epidemic Method ▾ Find in Topic

The Moving Epidemic Method



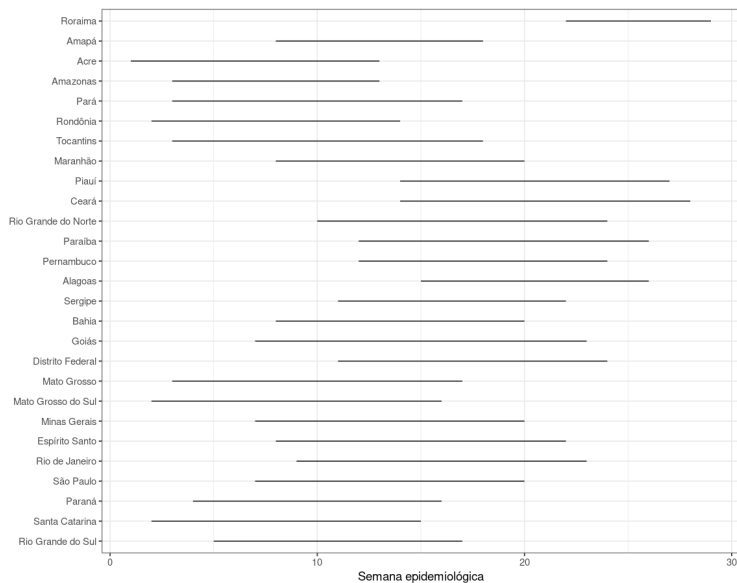
Documentation for package 'mem' version 2.16

- [DESCRIPTION file](#).

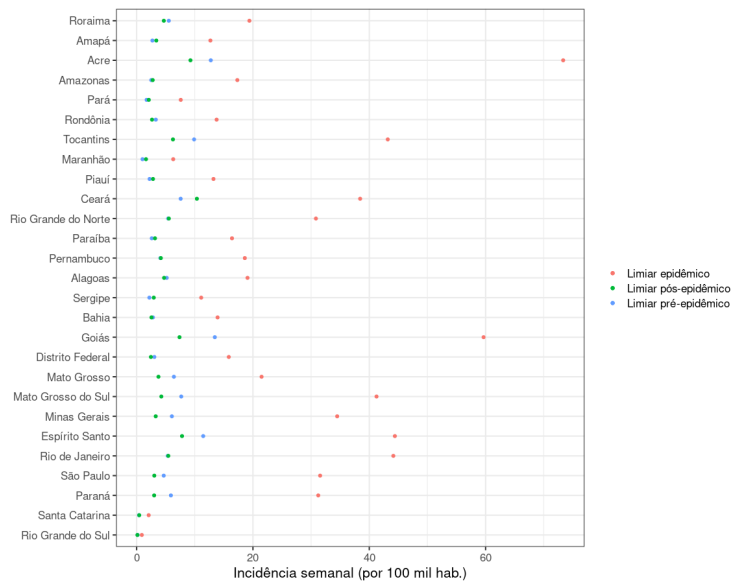
Help Pages

epimem	Deprecated function(s) in the mem package
eptiming	Deprecated function(s) in the mem package
flucyl	Castilla y Leon influenza crude rates
flucylraw	Castilla y Leon influenza standarised rates
full.series.graph	Creates the historical series graph of the datasets
memevolution	Evolution of estimators
memgoodness	Goodness of fit of the mem
memintensity	Thresholds for influenza intensity
memmodel	Methods for influenza modelization
memstability	Stability of indicators
memsurveillance	Creates the surveillance graph of the current season
memsurveillance.animated	Creates the animated graph of the surveillance of the current season
memtiming	Influenza Epidemic Timing
memtrend	Methods for influenza trend calculation
optimum.by.inspection	Inspection calculation of the optimum
processPlots	Full process plots for mem
roc.analysis	Analysis of different indicators to find the optimum value of the window parameter
summary.epidemic.plot.epidemic.print.epidemic	Influenza Epidemic Timing
summary.flu.plot.flu.print.flu	Methods for influenza modelization
transformdata	Data transformation
transformdata.back	Data transformation
transformseries	Transformation of series of data

Start and duration of an epidemic period (dengue by states in Brazil)



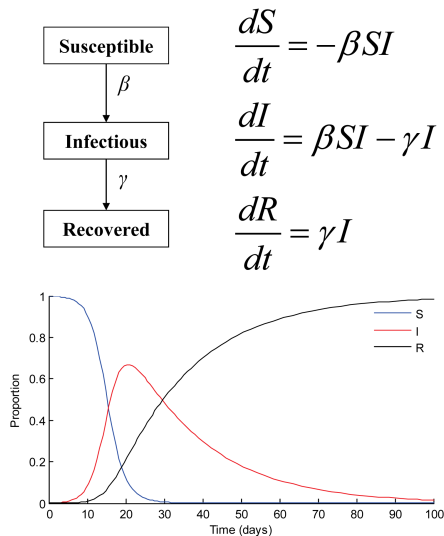
Dengue epidemic thresholds by states



- ① Introduction to infectious disease epidemiology
- ② Mathematical/deterministic models
 - SIR-like models
 - Semi-structured models
- ③ Real time analyses

A classical model

Luz, Struchiner & Galvani (2010).



Analysing a model

We have

$$S(t) + I(t) + R(t) = N \forall t$$

The basic reproductive number is

$$\mathcal{R}_0 = \frac{\beta N}{\gamma}. \quad (1)$$

Moreover,

$$\lim_{t \rightarrow \infty} I(t) = 0 \implies (S_e, 0, R_e),$$

is the only equilibrium point. Consider the Jacobian

$$J(S, I) = \begin{bmatrix} -\beta I & -\beta S \\ \beta I & \beta S - \gamma \end{bmatrix}.$$

Analysing a model (cont.)

- **Equilibria:** The characteristic polynomial is

$$\lambda^2 - (\beta S_e - \gamma) \lambda = 0,$$

thus $\lambda_1 = 0$ and $\lambda_2 = \beta S_e - \gamma$. This means we have neutral stability if $S_e = \gamma/\beta$, i.e., $\lambda_2 = 0$ and instability otherwise ($\lambda_2 > 0$ or $\lambda_2 < 0$).

- **Epidemic regimes:**

$$\frac{dI}{dS} = -1 + \frac{\gamma}{\beta} \frac{1}{S},$$

gives

$$I(t) = 1 - R(0) - S(t) + \frac{\gamma}{\beta} \ln \left(\frac{S(t)}{S(0)} \right).$$

From $\lim_{t \rightarrow \infty} I(t) = 0$ we know that

$$S(\infty) = 1 - R(0) + \frac{\gamma}{\beta} \ln \left(\frac{S(\infty)}{S(0)} \right).$$

Thus **an epidemic occurs iff**

$$\frac{\beta S(0)}{\gamma} > 1.$$

Big structured epidemic models

Coelho et al. (2020)

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

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

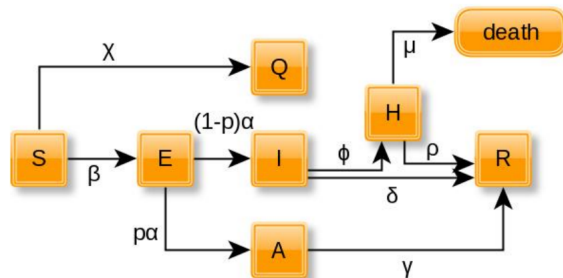
$$\frac{dI}{dt} = (1 - p)\alpha E - \delta I - \phi I,$$

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

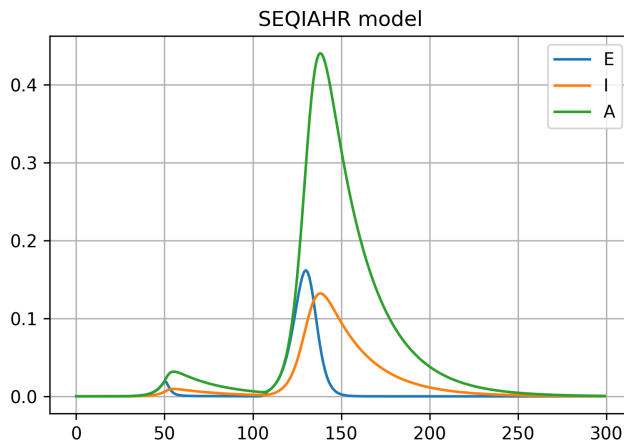
$$\frac{dH}{dt} = \phi I - (\rho + \mu)H,$$

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

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

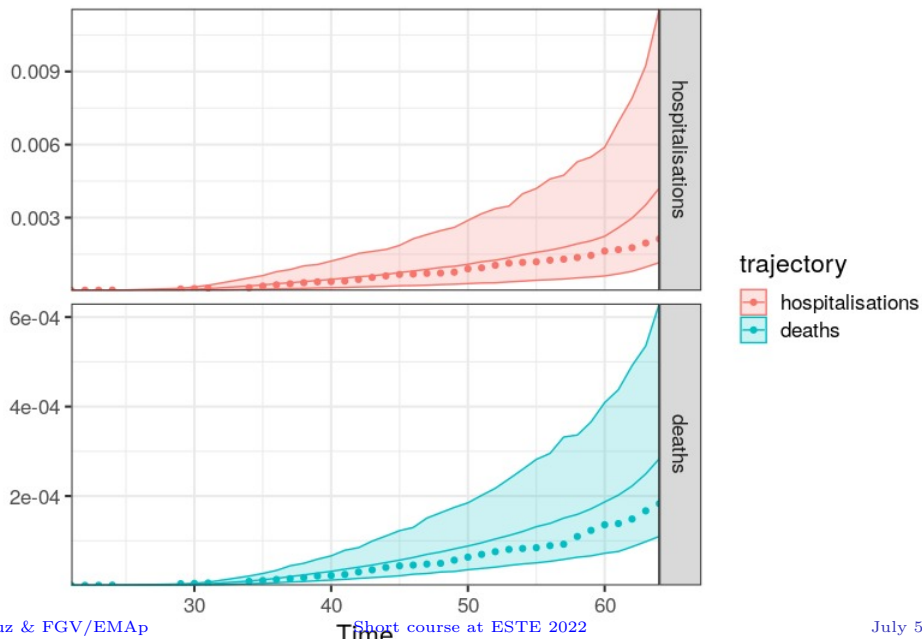


Big model (cont.)



$$\mathcal{R}_0 = \frac{\beta(1 - \xi)[p(\phi + \delta) + (1 - p)\gamma]}{\gamma(\delta + \phi)}$$

Why so rigid?



Stochastic models

Description	State transition	rate
Infection	$(S, E) \rightarrow (S - 1, E + 1)$	$\lambda(1 - \chi)S$
Exposed to I	$(E, I) \rightarrow (E - 1, I + 1)$	$(1 - p)\alpha E$
Exposed to A	$(E, A) \rightarrow (E - 1, A + 1)$	$p\alpha E$
Hospitalization	$(I, H) \rightarrow (I - 1, H + 1)$	ϕI
Recovery of I	$(I, R) \rightarrow (I - 1, R + 1)$	δI
Recovery of A	$(A, R) \rightarrow (A - 1, R + 1)$	γA
Recovery of H	$(H, R) \rightarrow (H - 1, R + 1)$	ρH
Death of H	$H \rightarrow H - 1$	μH

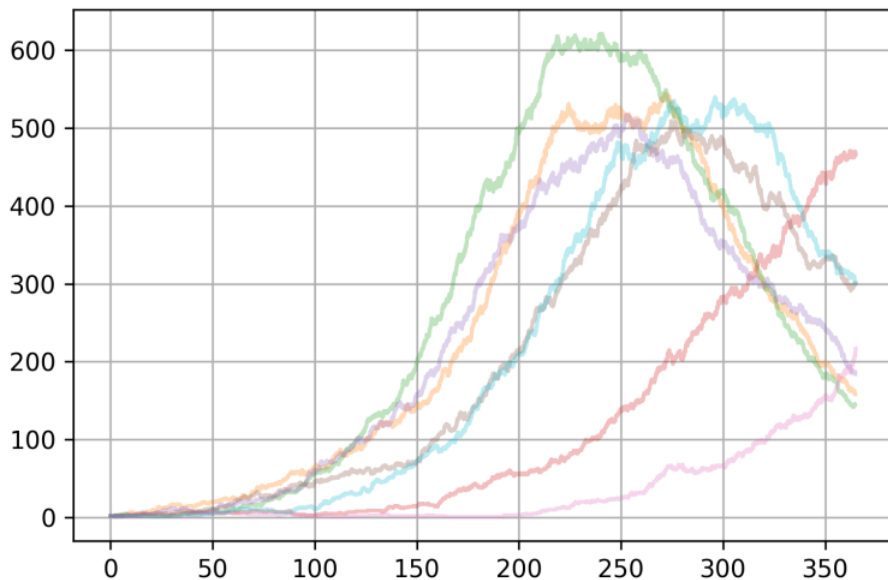
State probabilities

$$P_{s,e,i,a,h}(t) := \mathbb{P}(S = s, E = e, I = i, A = a, H = h).$$

FCK equation:

$$\begin{aligned} \frac{dP_{s,e,i,a,h}}{dt} = & P_{s+1,e-1,i,a,h}\lambda(1-\chi)(s+1) + P_{s,e+1,i-1,a,h}(1-p)\alpha(e+1) \\ & + P_{s,e+1,i,a-1,h}p\alpha(e+1) + P_{s,e,i+1,a,h-1}\phi(i+1) + P_{s,e,i+1,a,h}\delta(i+1) \\ & + P_{s,e,i,a+1,h}\gamma(a+1) + P_{s,e,i,a,h+1}\rho(h+1) + P_{s,e,i,a,h+1}\mu(h+1). \end{aligned}$$

Trajectories



Analysing a stochastic model

Probability-generating functions Starting with $I_i(0)$, the probability of an infected individual in state i producing offspring of type j given that $I_j(0)$ can be obtained from

$$f_i(z_1, \dots, z_k) = \sum_{j_k=0}^{\infty} \cdots \sum_{j_1=0}^{\infty} P_i(z_1, \dots, z_k) z_1^{j_1} \cdots z_k^{j_k}.$$

Now, define a matrix whose entries $m_{ji} = \frac{\partial f_i}{\partial u_i}|_{u=1}$ are the expected number of offspring generated in $i \rightarrow j$.

$$\mathbb{M} := \begin{bmatrix} 0 & 1-p & p \\ \frac{\beta(1-\chi)}{\beta(1-\chi)+\delta+\phi} & \frac{\beta(1-\chi)}{\beta(1-\chi)+\delta+\phi} & 0 \\ \frac{\beta(1-\chi)}{\beta(1-\chi)+\gamma} & 0 & \frac{\beta(1-\chi)}{\beta(1-\chi)+\gamma} \end{bmatrix}$$

Analysing a stochastic model (cont.)

Under some regularity conditions, we can calculate the **extinction probability**:

$$\mathbb{P}_0 = \prod_{i=1}^3 q_i^{k_i}$$

after finding $(q_1, q_2, q_3) \in (0, 1)^3$ which satisfy constraints. Here $k_1 = E(0)$, $k_2 = I(0)$ and $k_3 = A(0)$.

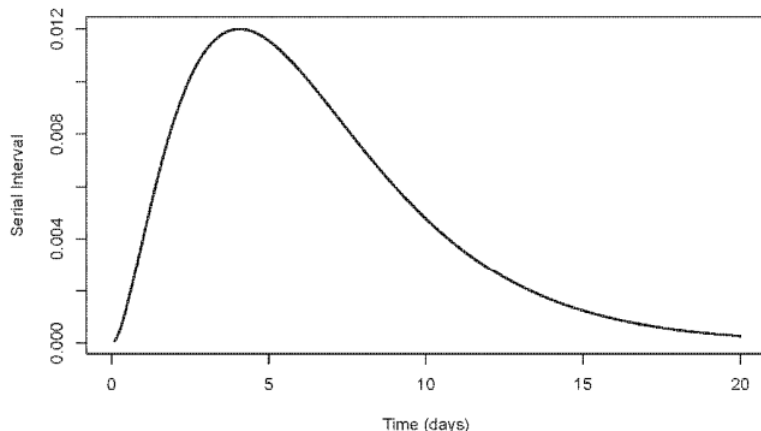
I_0	E_0	A_0	Approx. \mathbb{P}_0	SEIAHR \mathbb{P}_0	SIR \mathbb{P}_0
1	0	0	0.63	0.64	0.58
2	0	0	0.43	0.41	0.33
3	0	0	0.25	0.26	0.19
4	0	0	0.18	0.17	0.11
5	0	0	0.10	0.10	0.06

No more differential equations, please!

Start by looking at

$$R_t = \frac{I_t}{\sum_{s=1}^t I_{t-s} w_s}.$$

Then R_t is “the average number of secondary cases that each infected individual would infect if the conditions remained as they were at time t ” ([Cori et al. 2013](#)).



Modernising the model

$$R_{t,m} = R_{0,m} \left(2 \operatorname{logit}^{-1} \left(- \sum_{k=1}^4 (\alpha_k + \beta_{mk}) X_{ktm} \right) \right)$$

Priors:

$$\alpha_k \sim \operatorname{Normal}(0, 5);$$

$$\beta_{m,k} \sim \operatorname{Normal}(0, \gamma);$$

$$\gamma \sim \operatorname{HalfNormal}(0, 5);$$

$$R_{0,m} \sim \operatorname{Normal}(3.28, \kappa);$$

$$\kappa \sim \operatorname{HalfNormal}(0, 1/2).$$

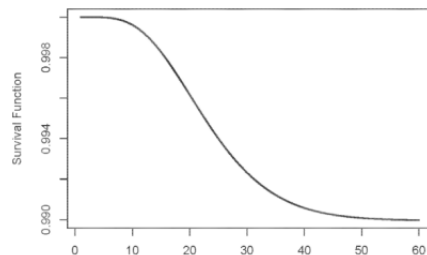
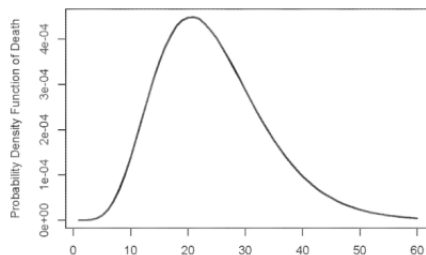
Modernising the model II: likelihood

$$c_{t,m} = R_{t,m} \sum_{\tau=0}^{t-1} c_{\tau,m} g_{t-\tau},$$

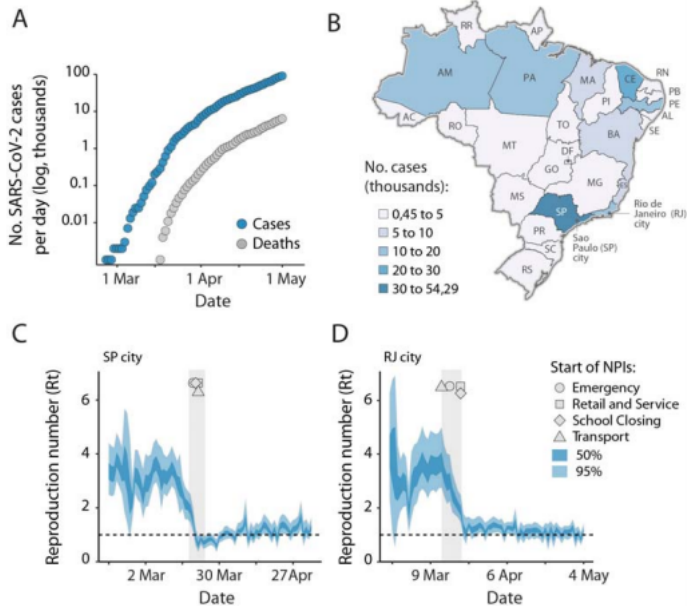
$$d_{t,m} = \sum_{\tau=0}^{t-1} c_{\tau,m} \pi_{t-\tau,m},$$

$$D_{t,m} \sim \text{NegativeBinomial}(d_{t,m}, \phi),$$

$$\phi \sim \text{HalfNormal}(0, 5).$$



The end product



- 1 Introduction to infectious disease epidemiology
- 2 Mathematical/deterministic models
- 3 Real time analyses
 - Nowcasting
 - Forecasting infectious diseases

Real time analysis


- In a (nearly) real time analysis of infectious disease cases, we aim to anticipate an outbreak/epidemic
- Usually this is part of an early warning system (EWS)
- An early warning system for outbreaks depends (from an stats point of view) on
 - Organised disease data systems
 - Periodic data release
 - Data cleaning and preparing
 - Monitoring system with reports, dashboards, etc to report the warning

Real time analysis

- The literature on EWS for chikungunya, dengue, malaria, yellow fever, and Zika outbreaks is scarce.

plos.org create account sign in

BROWSE PUBLISH ABOUT

SEARCH 


advanced search

PLOS NEGLECTED TROPICAL DISEASES

OPEN ACCESS


REVIEW

Early warning systems (EWSs) for chikungunya, dengue, malaria, yellow fever, and Zika outbreaks: What is the evidence? A scoping review

Laith Hussain-Alkhateeb , Tatiana Rivera Ramirez, Axel Kroeger, Ernesto Gozzer, Silvia Runge-Ranzinger

Published: September 16, 2021 • <https://doi.org/10.1371/journal.pntd.0009686>

41 Save	3 Citation
2,009 View	3 Share

Article	Authors	Metrics	Comments	Media Coverage
				

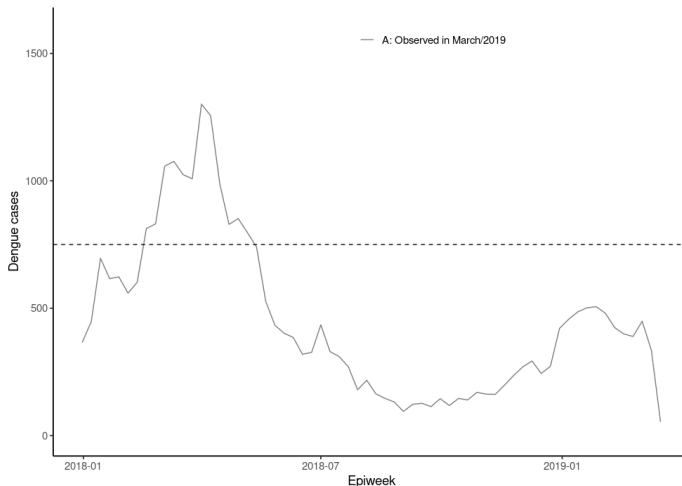
Download PDF

Print Share

Abstract

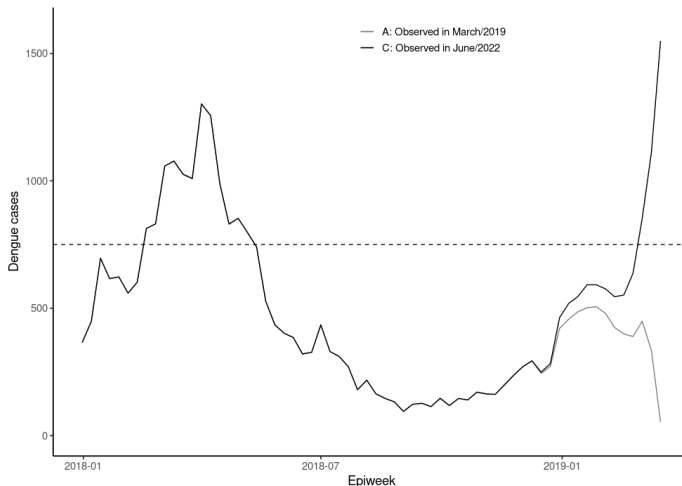
Dengue in Rio

Suppose today is 25/March/2019 and we estimate an epidemic threshold of 750 cases per week.



Dengue in Rio

Suppose today is 25/March/2019 and we estimate an epidemic threshold of 750 cases per week.



Notification delay



- Disease notification cases are in most of the cases delayed
 - Infrastructure issues
 - Epidemic level
 - Wrong diagnostics
 - etc

Data structure

Time	0	1	2	...	D-2	D-1	D	N	
1	$n_{1,0}$	$n_{1,1}$	$n_{1,2}$...	$n_{1,D-2}$	$n_{1,D-1}$	$n_{1,D}$	N_1	Observations
2	$n_{2,0}$	$n_{2,1}$	$n_{2,2}$...	$n_{2,D-2}$	$n_{2,D-1}$	$n_{2,D}$	N_2	
3	$n_{3,0}$	$n_{3,1}$	$n_{3,2}$...	$n_{3,D-2}$	$n_{3,D-1}$	$n_{3,D}$	N_3	
...	
T-D	$n_{T-D,0}$	$n_{T-D,1}$	$n_{T-D,2}$...	$n_{T-D,D-2}$	$n_{T-D,D-1}$	$n_{T-D,D}$	N_{T-D}	Nowcasting
T-D+1	$n_{T-D+1,0}$	$n_{T-D+1,1}$	$n_{T-D+1,2}$...	$n_{T-D+1,D-2}$	$n_{T-D+1,D-1}$	$n_{T-D+1,D}$	N_{T-D+1}	
T-D+2	$n_{T-D+2,0}$	$n_{T-D+2,1}$	$n_{T-D+2,2}$...	$n_{T-D+2,D-2}$	$n_{T-D+2,D-1}$	$n_{T-D+2,D}$	N_{T-D+2}	
T-2	$n_{T-2,0}$	$n_{T-2,1}$	$n_{T-2,2}$...	$n_{T-2,D-2}$	$n_{T-2,D-1}$	$n_{T-2,D}$	N_{T-2}	Forecasting
T-1	$n_{T-1,0}$	$n_{T-1,1}$	$n_{T-1,2}$...	$n_{T-1,D-2}$	$n_{T-1,D-1}$	$n_{T-1,D}$	N_{T-1}	
T	$n_{T,0}$	$n_{T,1}$	$n_{T,2}$...	$n_{T,D-2}$	$n_{T,D-1}$	$n_{T,D}$	N_T	
T+1	$n_{T+1,0}$	$n_{T+1,1}$	$n_{T+1,2}$...	$n_{T+1,D-2}$	$n_{T+1,D-1}$	$n_{T+1,D}$	N_{T+1}	Forecasting
T+2	$n_{T+2,0}$	$n_{T+2,1}$	$n_{T+2,2}$...	$n_{T+2,D-2}$	$n_{T+2,D-1}$	$n_{T+2,D}$	N_{T+2}	
...	
T+K	$n_{T+K,0}$	$n_{T+K,1}$	$n_{T+K,2}$...	$n_{T+K,D-2}$	$n_{T+K,D-1}$	$n_{T+K,D}$	N_{T+K}	

- Total number of cases per week is given by the sum of columns

The chain-ladder model

- The likelihood

$$\log(n_{t,d}) \sim N(m_{t,d}, \tau),$$

$$m_{t,d} = \mu + \alpha_t + \beta_d$$

for $t = 1, 2, \dots, n$, $d = 0, 1, \dots, D$.

- It is a two-way ANOVA, proposed by Renshaw (1989) with a Bayesian version proposed by Verral (1989a)
- Verral (1989b) propose a state-space representation
- Mack (1993) proposed a Poisson/glm approach with fixed effects

Correcting delays

Using historical data, reporting delays can be corrected

- Bastos et al. (2016, 2019) adapted the chain-ladder model adding spatiotemporal random effects (INLA)

Correcting delays

Using historical data, reporting delays can be corrected

- Bastos et al. (2016, 2019) adapted the chain-ladder model adding spatiotemporal random effects (INLA)
- McGough et al. (2020) also adapted the chain-ladder model with temporal random effects (JAGS + R package NoBS)

Correcting delays

Using historical data, reporting delays can be corrected

- Bastos et al. (2016, 2019) adapted the chain-ladder model adding spatiotemporal random effects (INLA)
- McGough et al. (2020) also adapted the chain-ladder model with temporal random effects (JAGS + R package NoBS)
- Rotejanaprasert et al. (2020) spatiotemporal sliding windows (INLA)

Correcting delays

Using historical data, reporting delays can be corrected

- Bastos et al. (2016, 2019) adapted the chain-ladder model adding spatiotemporal random effects (INLA)
- McGough et al. (2020) also adapted the chain-ladder model with temporal random effects (JAGS + R package NoBS)
- Rotejanaprasert et al. (2020) spatiotemporal sliding windows (INLA)
- Miller et al. (2022) use google searches to improve the correction (INLA)

Bastos et al. (2019) model

- The likelihood

$$n_{t,d} \sim \text{NegBin}(\lambda_{t,d}, \phi), \quad \lambda_{t,d} > 0, \quad \phi > 0.$$

Bastos et al. (2019) model

- The likelihood

$$n_{t,d} \sim \text{NegBin}(\lambda_{t,d}, \phi), \quad \lambda_{t,d} > 0, \quad \phi > 0.$$

- Adding fixed and random effects

$$\log(\lambda_{t,d}) = \mu + \alpha_t + \beta_d + \mathbf{x}'_{t,d}\boldsymbol{\gamma},$$

for $t = 1, 2, \dots, n$, $d = 0, 1, \dots, D$, and $\mathbf{x}_{t,d}$ a vector covariates.

Bastos et al. (2019) model

- The likelihood

$$n_{t,d} \sim \text{NegBin}(\lambda_{t,d}, \phi), \quad \lambda_{t,d} > 0, \quad \phi > 0.$$

- Adding fixed and random effects

$$\log(\lambda_{t,d}) = \mu + \alpha_t + \beta_d + \mathbf{x}'_{t,d}\boldsymbol{\gamma},$$

for $t = 1, 2, \dots, n$, $d = 0, 1, \dots, D$, and $\mathbf{x}_{t,d}$ a vector covariates.

- Time (α_t) and delay (β_d) random effects may be modelled using
 - First or second order Gaussian random walks (RW1, RW2)
 - Gaussian autoregressive processes (AR(p))
- Finally, weakly informative priors are set for $(\mu, \phi, \boldsymbol{\gamma}, \boldsymbol{\psi}_\alpha, \boldsymbol{\psi}_\beta, \dots)$

Inference (Nowcast)

- Learn about all parameters throughout the posterior distribution
- Access the predictive distribution of unknown $n_{t,d}$ s (The run-off triangle)

$$p(n_{t^*,d^*} \mid \{n_{t,d}, t+d < T\}), \quad \{(t^*, d^*) : T < t^* + d^* < T + D\}.$$

- Derive the predictive distribution of the total notifications at time t

$$p(N_t = \sum_d n_{t,d} \mid \{n_{t,d}, t+d < T\}).$$

Inference (Monte Carlo)

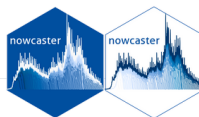
Sample from the predictive distribution as the following:

- 1 Sample $(\phi, \mu, \alpha_t, \beta_d)$ from the joint posterior;
- 2 Sample an unknown $n_{t,d}$ from the likelihood using the sampled parameters;
- 3 Calculate $N_t = \sum_d n_{t,d}$ (some values are known, others have just been sampled)

Implemented in the R package *nowcaster*

nowcaster

Nowcaster



nowcaster is a R package for “nowcasting” epidemiological time-series. Every single system of notification has an intrinsic delay, **nowcaster** can estimate how many counts of any epidemiological data of interest (*i.e.*, daily cases and deaths counts) by fitting a negative binomial model to the time steps of delay between onset date of the event, (*i.e.*, date of first symptoms for cases or date of occurrence of death) and the date of report (*i.e.*, date of notification of the case or death).

nowcaster is based on the [R-INLA](#) and [INLA](#) packages for “Integrated Nested Laplace Approximation” algorithm to Bayesian inference. **INLA** is a fast alternative to others methods for Bayesian inference like **MCMC**. An introduction to **INLA** can be found [here](#).

nowcaster is build for epidemiological emergency use, it was constructed for the Brazilian Severe Acute Respiratory Illness (SARI) surveillance database (SIVEP-Gripe).

Installing

Before installing the package certify you have an active installation of **INLA**, to do so you can run the follwing code:

```
install.packages("INLA",
  repos=c(getOption("repos"),
    INLA="https://inla.r-inla-download.org/R/stable"),
  dep=TRUE)
```

<https://covid19br.github.io/nowcaster/>

Links

[Browse source code](#)

[Report a bug](#)

License

[Full license](#)


GPL (>= 3)

Citation

[Citing nowcaster](#)

Developers

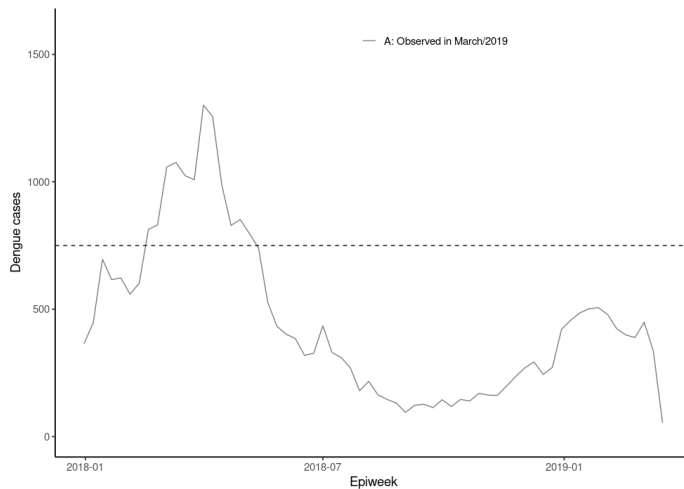
[Rafael Lopes](#)

Author, maintainer 

[Leonardo Bastos](#)

Author 

Dengue in Rio



In R

- Data structure (input)

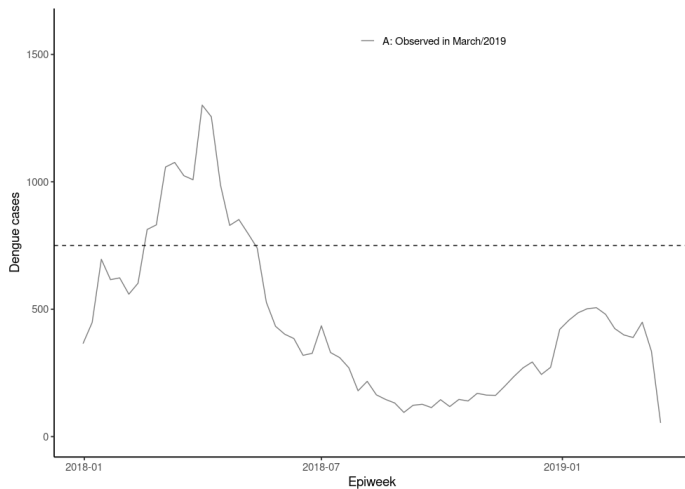
```
> head(dengue.RJ)
  SG_UF ID_MN_RESI DT_DIGITA DT_SIN_PRI DT_NOTIFIC CS_SEXO NU_IDADE_N
1   33    330270 2019-01-10 2018-08-25 2018-09-05      F      4062
2   33    330270 2019-01-08 2018-12-12 2018-12-19      F      4053
3   33    330270 2019-01-10 2018-06-23 2018-06-29      F      4058
4   33    330270 2018-11-27 2018-06-12 2018-06-15      F      4012
5   33    330270 2018-11-14 2018-05-25 2018-06-04      M      4008
6   33    330270 2018-11-12 2018-07-01 2018-07-11      F      4058
```

- running the nowcasting model

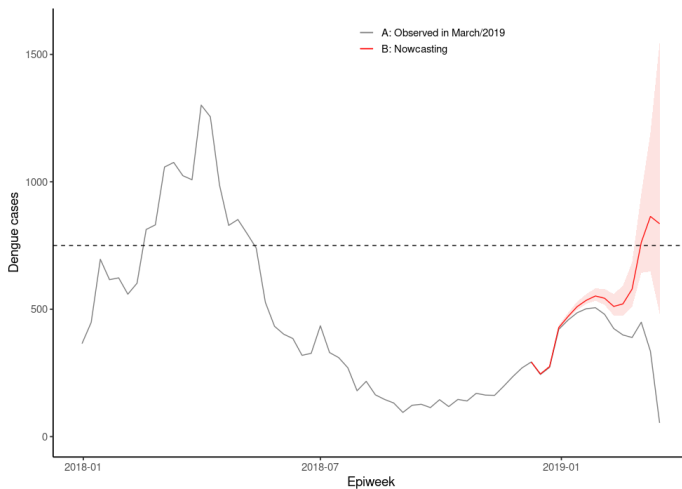
```
library(nowcaster)

dengue.RJ.now <- nowcasting_inla(dataset = dengue.RJ %>%
  filter(DT_DIGITA < "2019-03-25"),
  date_onset = DT_SIN_PRI,
  date_report = DT_DIGITA,
  data.by.week = T )
```

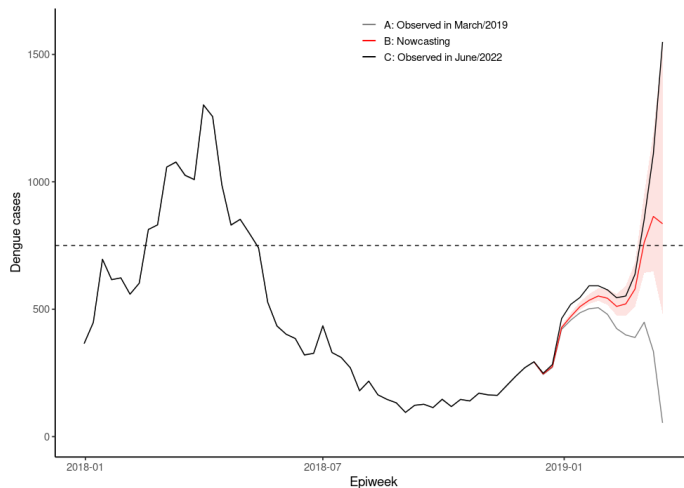
Dengue in Rio



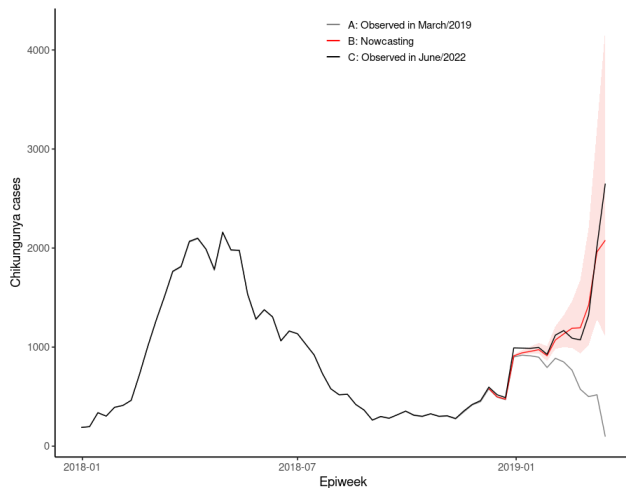
Dengue in Rio



Dengue in Rio



Chikungunya in Rio



Infodengue

- Infodengue is an early warning system for dengue, chikungunya and Zika



Início

Saiba Mais

Equipe

Participe

Dados

Relatórios

Encontre um município ...

Login

Situação de casos estimados

Análise integrada de dados epidemiológicos, climáticos e redes sociais.

Atualização semanal:

- Incidência estimada (nowcasting)
- Cidades com condições favoráveis para transmissão
- Cidades em níveis de atenção



Dengue
SE 41-44/2021

Funcionalidades:

- Relatórios municipais
- Mapas estaduais
- API

Participe:

Existem várias formas de participar,
[Confira aqui!](#)

<https://info.dengue.mat.br>

Infodengue

- (Input) (Nearly) real time data
 - Notified cases (instead of suspected/confirmed cases)
 - Social network data (tweets)
 - Real time weather data

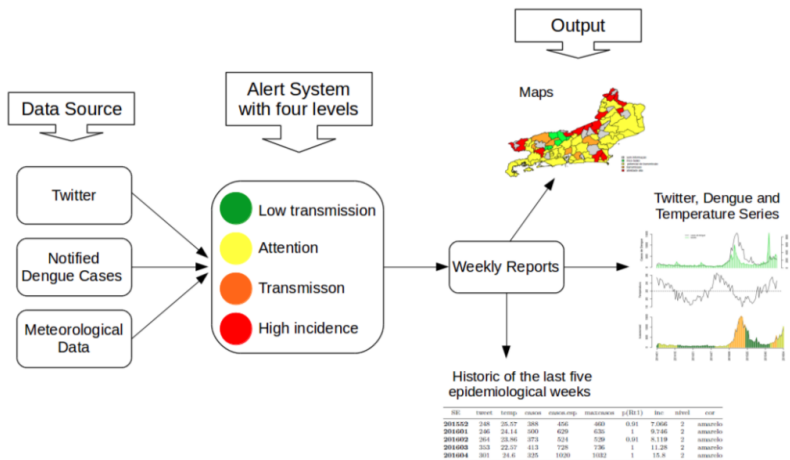
Infodengue

- (Input) (Nearly) real time data
 - Notified cases (instead of suspected/confirmed cases)
 - Social network data (tweets)
 - Real time weather data
- (Output)
 - Reports at city level
 - Corrected estimates
 - Four-color early warning mechanism

Infodengue

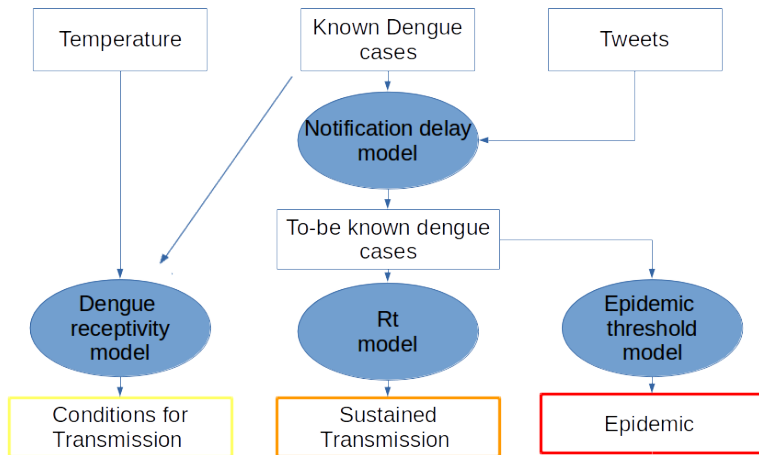
- (Input) (Nearly) real time data
 - Notified cases (instead of suspected/confirmed cases)
 - Social network data (tweets)
 - Real time weather data
- (Output)
 - Reports at city level
 - Corrected estimates
 - Four-color early warning mechanism
- (How?)
 - Data cleaning
 - Defining epidemic thresholds
 - Delay correction

Infodengue: Data workflow



<https://info.dengue.mat.br>

Infodengue: Model workflow



<https://info.dengue.mat.br>

Infodengue: Issues

- Big data: Weekly analysis for all 5500+ Brazilian municipalities

Infodengue: Issues

- Big data: Weekly analysis for all 5500+ Brazilian municipalities
- Low frequency cases versus high frequency climate data

Infodengue: Issues

- Big data: Weekly analysis for all 5500+ Brazilian municipalities
- Low frequency cases versus high frequency climate data
- Multivariate modelling (three diseases transmitted by the same mosquito)

Infodengue: Issues

- Big data: Weekly analysis for all 5500+ Brazilian municipalities
- Low frequency cases versus high frequency climate data
- Multivariate modelling (three diseases transmitted by the same mosquito)
- Misclassification (diseases with similar symptoms)

Infodengue: Issues

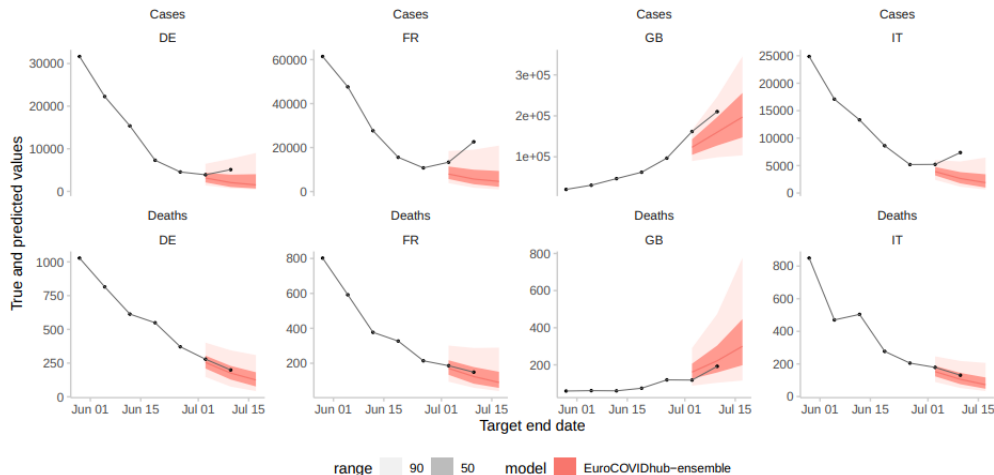
- Big data: Weekly analysis for all 5500+ Brazilian municipalities
- Low frequency cases versus high frequency climate data
- Multivariate modelling (three diseases transmitted by the same mosquito)
- Misclassification (diseases with similar symptoms)
- (Good) forecasting models (short 4-week ahead and long term one-year ahead)

- ① Introduction to infectious disease epidemiology
- ② Mathematical/deterministic models
- ③ Real time analyses
 - Nowcasting
 - Forecasting infectious diseases

How do you assess predictions from a model?

Desiderata:

- (Well-calibrated) Probabilistic predictions;
- Encourage careful and honest predictions;



Fixing notation

Presentation will be mostly based on [Gneiting & Raftery \(2007\)](#) and [Bosse et al. \(2022\)](#). Let \mathcal{P} be a convex class of probability measures on (Ω, \mathcal{A}) . We call $P \in \mathcal{P}$ a probabilistic forecast.

Definition (Scoring rule)

We say $S(P, \cdot) : \Omega \rightarrow [-\infty, \infty]$ is a **scoring rule** if it is measurable and $S(P, \cdot)$ is P -quasi-integrable for all $P \in \mathcal{P}$.

The expected score under $Q \in \mathcal{P}$ if the forecast is P is

$$S(P, Q) := \int_{\Omega} S(P, \omega) dQ(\omega).$$

Definition (Strictly proper scoring rule)

We say S is **proper** if $S(Q, Q) \geq S(P, Q)$ for all $P, Q \in \mathcal{P}$. In addition, we say S is **strictly proper** if equality is achieved only for $P = Q$.

Categorical examples

- **Spherical score.**

For $\alpha > 1$ we can define

$$S(\mathbf{p}, i) = \frac{p_i^{\alpha-1}}{\left(\sum_{j=1}^m p_j^\alpha\right)^{\frac{\alpha-1}{\alpha}}} \quad (2)$$

- **Logarithmic score.**

When $G(\mathbf{p})$ is the Shannon entropy, we have $S(\mathbf{p}, i) = \log p_i$ and

$$d(\mathbf{p}, \mathbf{q}) = \sum_{j=1}^m \log \left(\frac{q_j}{p_j} \right), \quad (3)$$

as the Kullback-Leibler divergence.

Scoring rules for density forecasts

Let μ be a σ -finite measure on (Ω, \mathcal{A}) . For $\alpha > 1$ define \mathcal{L}_α be the space of probability measures on (Ω, \mathcal{A}) such that $\nu \ll \mu$ and $p(\omega) = \frac{d\nu}{d\mu}(\omega)$ and

$$\|p\|_\alpha = \left(\int_\Omega p(\omega)^\alpha d\mu(\omega) \right)^\alpha < \infty.$$

We establish a correspondence between the forecast P and its μ -density, p . **Examples:**

- **Quadratic:**

$$\text{QS}(p, \omega) = 2p(\omega) - \|p\|_2^2, \quad (4)$$

is strictly proper relative to \mathcal{L}_2 class of probability measures.

- **Pseudo-spherical:**

$$\text{PseudoS}(p, \omega) = \frac{p(\omega)^{\alpha-1}}{\|p\|_\alpha^{\alpha-1}}, \quad (5)$$

- **Logarithmic score:**

$$\text{LogS}(p, \omega) = \log p(\omega), \quad (6)$$

is what happens to the pseudo-spherical score when $\alpha \rightarrow 1$.

Continuous ranked probability score

The continuous ranked probability score (CRPS):

$$\text{CRPS}(F, x) = - \int_{-\infty}^{\infty} (F(y) - \mathbb{I}\{y \geq x\})^2 dy, \quad (7)$$

can be seen as the integral of the Brier scores for the associated binarisation of the forecasts based on x as cutoff.

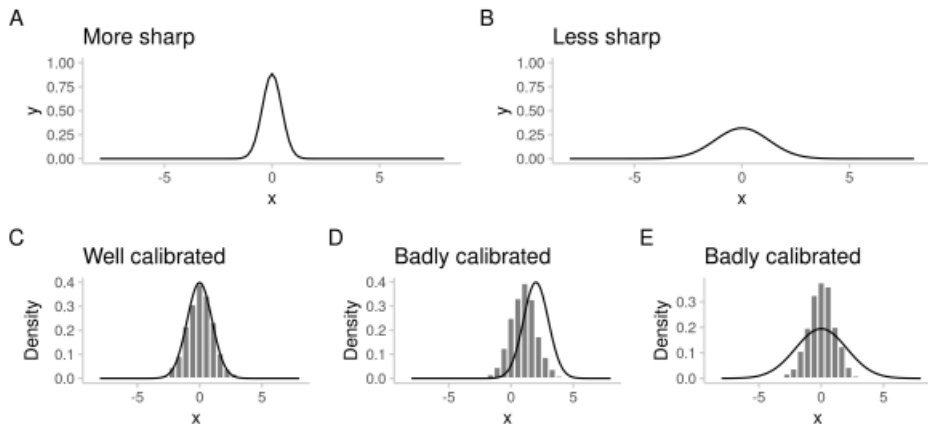
Example:

$$\text{CRPS}(\text{Normal}(\mu, \sigma^2), x) = \sigma \left[\frac{1}{\sqrt{\pi}} - 2\varphi\left(\frac{x - \mu}{\sigma}\right) - \frac{x - \mu}{\sigma} \left(2\Phi\left(\frac{x - \mu}{\sigma}\right) - 1 \right) \right],$$

where φ and Φ are the probability density function and cumulative distribution function of a standard normal, respectively.

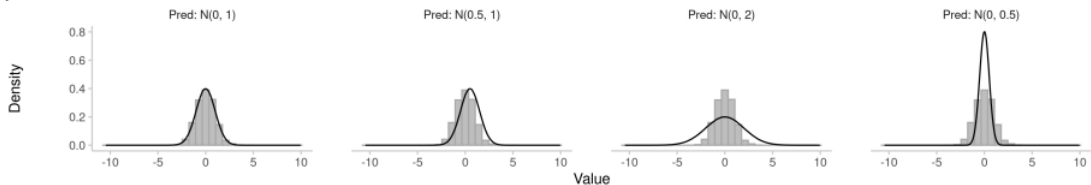
Being right and being sure of it

- Calibration
- Sharpness

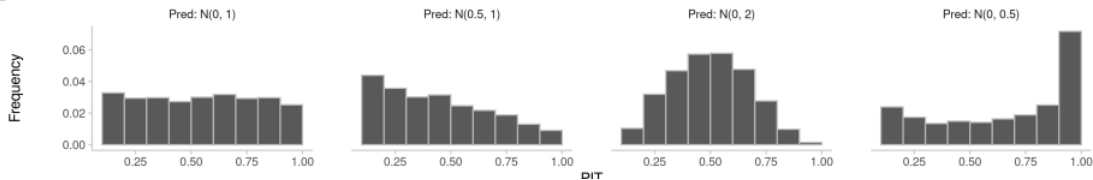


True target is Normal(0, 1)

A

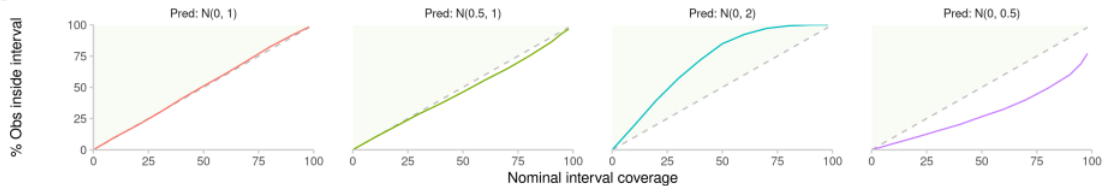


B

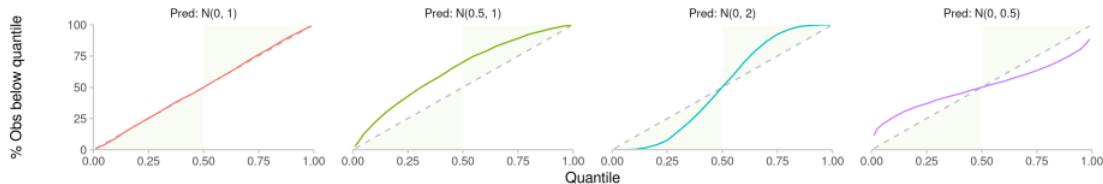


True target is Normal(0, 1)

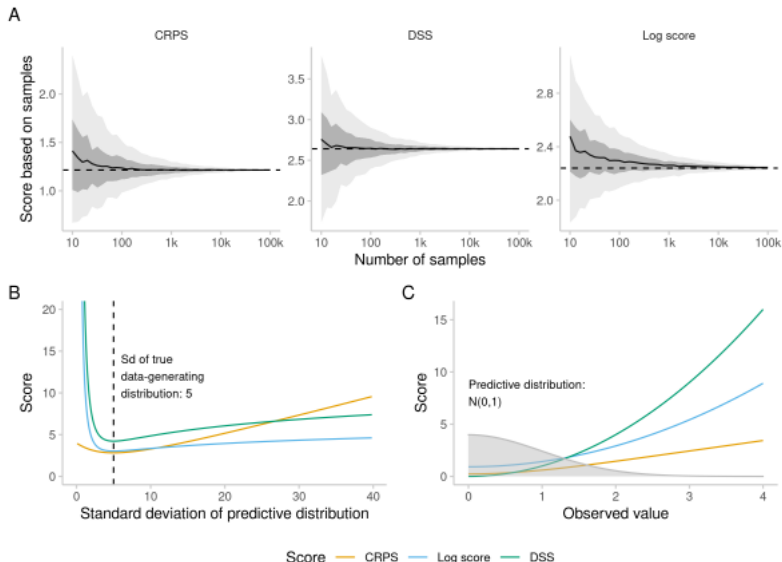
C



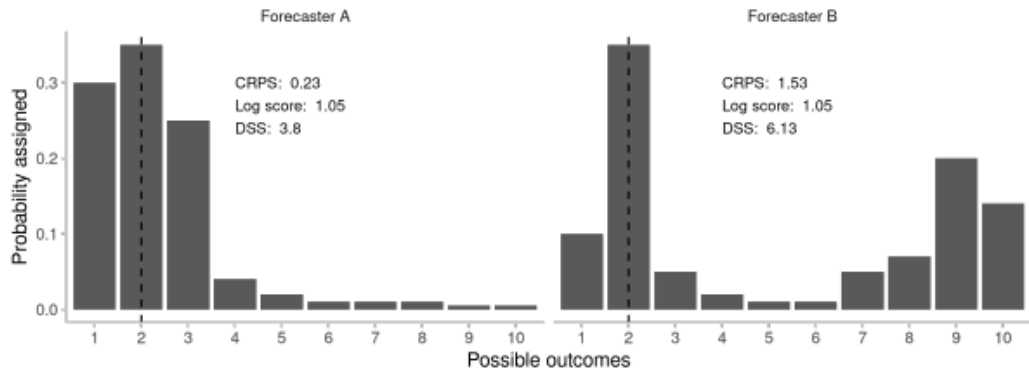
D



Which score to pick if you're predicting hospitalisations?



Local *vs* Global



Scale dependence

