# Statistical modelling of climate-sensitive diseases

**Ania Kawiecki Peralta (ania.kawiecki@bsc.es)**

**Carles Milà Garcia (carles.milagarcia@bsc.es)**

Global Health Resilience group

*Advanced Webinar Series on Spatiotemporal Modeling of Climate-Sensitive Diseases*
*28st of January 2026*

Earth Sciences Department | Barcelona Supercomputing Center Centro Nacional de Supercomputación

Global Health Resilience

**Ania Kawiecki Peralta**



I am a postdoc at the Global Health Resilience group at the BSC. My background is in Veterinary Medicine, and I have a PhD in Epidemiology studying dengue virus vector surveillance and control.

Currently I am working on developing R packages to facilitate disease risk modeling and prediction using Bayesian spatio-temporal models in INLA.

**Carles Milà**



I am a data scientist at the Global Health Resilience group at the BSC. My background is in statistics and geoinformatics, and I have a PhD in spatial modelling for exposure assessment.

I am currently working on developing R packages for climate-sensitive data processing and modelling.

1. Introduction: Linear model recap

2. Hierarchical generalized linear models (from a Bayesian perspective)
   - Introduction to generalized linear models
   - Basics of Bayesian inference
   - Hierarchical models

3. Model terms in a spatiotemporal context

4. Forecasting for early warning systems

5. Questions at the end!

# Linear model recap

# Linear model revisited

```
disease_cases ~ rainfall + mean_temperature
```
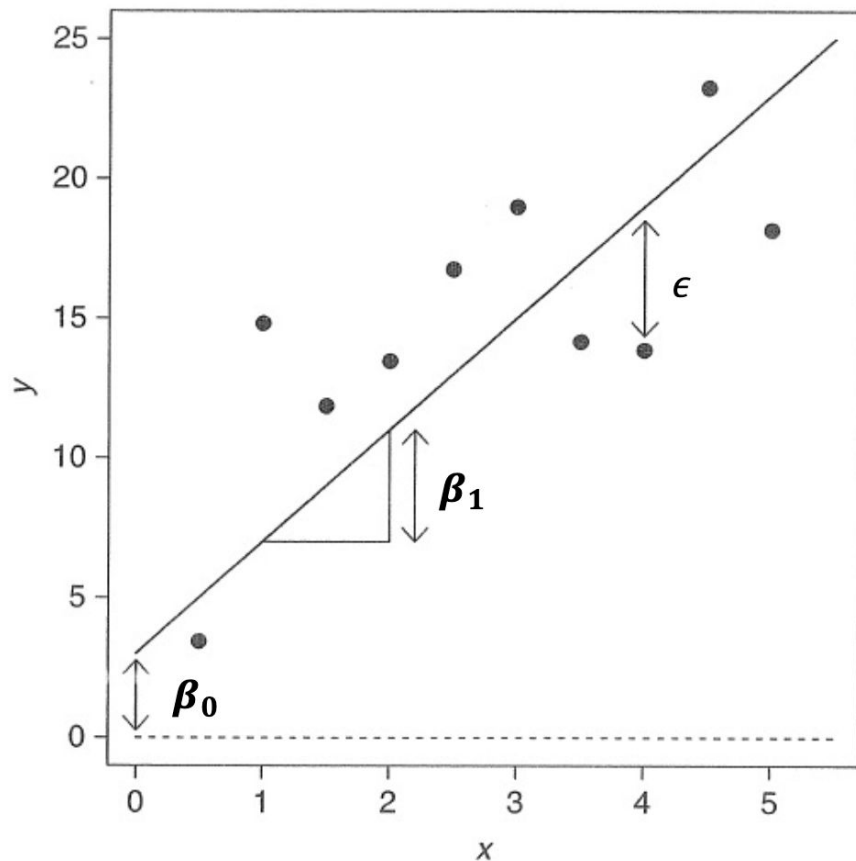
$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \varepsilon_i$$

$\beta_0$ (intercept): number of dengue cases expected if rainfall = 0 and mean temperature = 0

$\beta_1$ = how many more cases you expect per 1 unit increase in rainfall, holding mean temperature constant

$\beta_2$ = how many more cases you expect per 1 unit increase in mean temperature, holding rainfall constant

$\varepsilon_i$ = error term

# Linear model revisited

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \varepsilon_i$$

where

$$\varepsilon_i \sim N(0, \sigma)$$

What are (some) of the assumptions in this model?

- Linearity in the predictors: In the last webinar we saw that some variables have non-linear effects.

- Conditional independence of the observations: Our data are structured in space and time and therefore have autocorrelation.

- The response is Normally distributed conditional on the predictors and parameters: Not true for disease case counts.

$$Y_i | \mu_i, \sigma^2 \sim N(\mu_i, \sigma^2)$$
$$\mu_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2}$$

$\mu$: mean of the distribution

$\sigma^2$: variance of the distribution

```
disease_cases ~ rainfall + mean_temperature
```

## There's room for improvement!

| | | |
|---|---|---|
| How can we take into account the disease **seasonality**? | How can we account for differences and correlation between **spatial areas**? | Can we improve the model to reflect the **distribution** of the case count data? |
| How can we take into account the **interannual variation** in cases? | How can we incorporate **non-linear** relationships? | How do we use this model for **forecasting** and to account for **uncertainty**? |

# Hierarchical generalized linear models (from a Bayesian perspective)

- ○ Introduction to generalized linear models

- ○ Basics of Bayesian inference

- ○ Hierarchical models

Our linear model:

$$Y_i | \mu_i, \sigma^2 \sim N(\mu_i, \sigma^2)$$
$$\mu_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2}$$

Some reasons this is problematic to predict disease case counts:

- We could predict negative cases (e.g. -2)

- We could predict non-integer cases (e.g. 4.3)

- Assumes that the variance is constant and is independent from the mean.

  - We know that the variance of a count scales with the mean: If the mean count $\mu$ is large, the variance $\sigma^2$ will also be larger.

Enter the Generalized Linear Model (GLM)

$\theta$: Other parameters of the model

$f$: A distribution in the *Exponential family*

$\mu$: The mean of the distribution
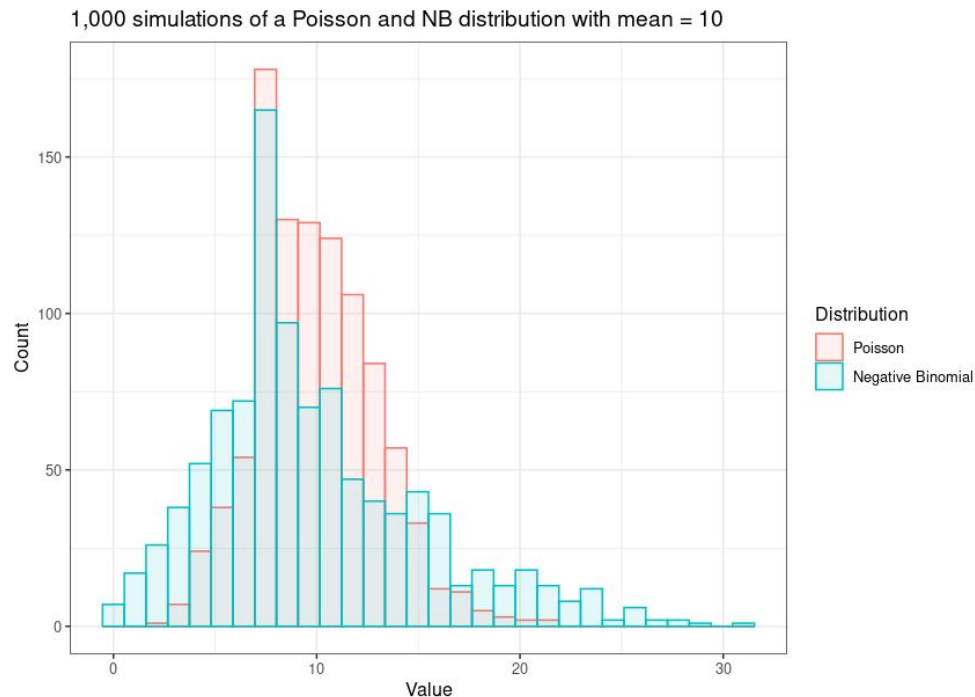
$$Y_i | \mu_i, \theta \sim f(\mu_i, \theta)$$

$$g(\mu_i) = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2}$$

$g$: Link function

Linear predictor

Distributions suitable for count data (integer support):

- **Poisson:**
  variance = mean

- **Negative Binomial:**
  variance > mean
  (i.e. *overdispersion*)



1,000 simulations of a Poisson and NB distribution with mean = 10

Distribution
- Poisson
- Negative Binomial

Disease case counts usually exhibit *overdispersion*:
Negative Binomial is often used

As a link function, we use the *log*. Therefore, our GLM tailored for case counts becomes:

$$Y_i | \mu_i, \theta \sim \text{NegBin}(\mu_i, \theta)$$

$$log(\mu_i) = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2}$$

Why use the *log* as link function?

- Ensures that the mean is positive.

$$\mu_i = exp(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2})$$

- Multiplicative effect of covariates, useful to capture skewed case counts.

$$\mu_i = exp(\beta_0) \cdot exp(\beta_1 x_{i1}) \cdot exp(\beta_2 x_{i2})$$
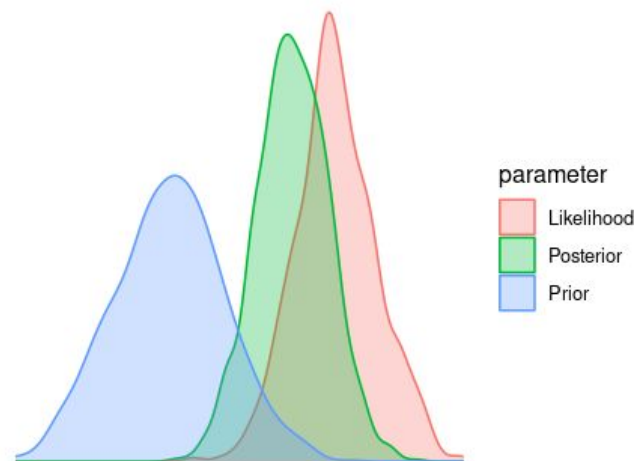
How can we interpret the model coefficients $\square_1$ and $\square_2$ ?

$$Y_i | \mu_i, \theta \sim \text{NegBin}(\mu_i, \theta)$$
$$log(\mu_i) = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2}$$

- $\square_1$: increasing the temperature by 1 makes the $\log(\mu_i)$ increase by $\square_1$

  - $\square_1 < 0$ means a decrease in risk while $\square_1 > 0$ means an increase in risk.

- $\exp(\square_1)$: can be interpreted as the multiplicative factor on the mean count per unit increase in temperature. Why?

$$\mu_i(x_{i1} + 1) = \exp\left(\beta_0 + \beta_1(x_{i1} + 1) + \beta_2 x_{i2}\right)$$
$$= \exp\left(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2}\right) \cdot \exp(\beta_1)$$
$$= \mu_i(x_{i1}) \cdot \exp(\beta_1)$$

We have disease case counts that vary in space and time:

- The population at risk varies in time and space so it's difficult to compare counts.
- Could we standardize them somehow?

```
Model counts  →  Model rates
```

Population at risk

$$Y_i | \mu_i, \theta \sim \text{NegBin}(\mu_i, \theta)$$

$$log(\frac{\mu_i}{P_i}) = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2}$$

$$log(\mu_i) = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + log(P_i)$$

We add a population offset to model rates rather than counts

Bayesian inference allows us to estimate model parameters while characterizing their uncertainty through their **full probability distributions**.

$$Y_i | \mu_i, \theta \sim \text{NegBin}(\mu_i, \theta)$$
$$log(\mu_i) = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + log(P_i)$$

3 main components in Bayesian estimation:

- Prior distribution: Our **prior beliefs** about the parameters before observing the data.

- Likelihood: How probable the **observed data** are given the model parameters.

- Posterior distribution: **Updated beliefs** about the parameter after observing the data.

Bayesian inference updates the prior using the information in the data (likelihood) to get the posterior.



parameter
Likelihood
Posterior
Prior

- Priors are specified **before fitting the model** to the current data

- Unless we have very strong evidence, <u>weakly informative priors</u> are often a good default.

- Weakly informative priors:

  - Allow the data to dominate when information is strong.



Weakly informative prior for □

# Bayesian inference : posterior predictive distribution

Why **Bayesian** in early warning systems?

We also obtain a probability distribution for the predictions: **the posterior predictive distribution.**



We can calculate uncertainty intervals directly in the distribution: *credible intervals*

We can calculate the probability of exceeding a threshold: *outbreak probability*

# Hierarchical models

| date | micro_name | dengue_cases | tmin | pdsi | meso_name | water_network |
|---|---|---|---|---|---|---|
| 2001–01–01 | Alto Taquari | 0 | 22.33064 | –0.67288578 | Centro Norte De Mato Grosso Do Sul | 86.21000 |
| 2001–02–01 | Alto Taquari | 2 | 22.09503 | –0.79167610 | Centro Norte De Mato Grosso Do Sul | 86.21000 |
| 2001–03–01 | Alto Taquari | 6 | 21.65975 | 0.19557676 | Centro Norte De Mato Grosso Do Sul | 86.21000 |
| ⋮ | | | | | | |
| 2001–01–01 | Aquidauana | 1 | 22.94171 | 1.60573995 | Pantanais Sul Mato–Grossense | 84.18500 |
| 2001–02–01 | Aquidauana | 0 | 22.75295 | 2.24432206 | Pantanais Sul Mato–Grossense | 84.18500 |
| 2001–03–01 | Aquidauana | 5 | 22.04041 | 1.54781866 | Pantanais Sul Mato–Grossense | 84.18500 |
| ⋮ | | | | | | |
| 2001–01–01 | Baixo Pantanal | 0 | 23.50009 | 0.08895861 | Pantanais Sul Mato–Grossense | 84.18500 |
| 2001–02–01 | Baixo Pantanal | 1 | 23.27970 | 0.25999102 | Pantanais Sul Mato–Grossense | 84.18500 |
| 2001–03–01 | Baixo Pantanal | 1 | 22.71204 | 0.26609710 | Pantanais Sul Mato–Grossense | 84.18500 |

# Hierarchical models

*i*: individual-level

*j*: group-level

| date | micro_name | dengue_cases | tmin | pdsi | meso_name |
|------|-----------|-------------|------|------|-----------|
| 2001-01-01 | Alto Taquari | 0 | 22.33064 | -0.67288578 | Centro Norte De Mato Grosso Do Sul |
| 2001-02-01 | Alto Taquari | 2 | 22.09503 | -0.79167610 | Centro Norte De Mato Grosso Do Sul |
| 2001-03-01 | Alto Taquari | 6 | 21.65975 | 0.19557676 | Centro Norte De Mato Grosso Do Sul |
| ⋮ | | | | | |
| 2001-01-01 | Aquidauana | 1 | 22.94171 | 1.60573995 | Pantanais Sul Mato-Grossense |
| 2001-02-01 | Aquidauana | 0 | 22.75295 | 2.24432206 | Pantanais Sul Mato-Grossense |
| 2001-03-01 | Aquidauana | 5 | 22.04041 | 1.54781866 | Pantanais Sul Mato-Grossense |
| ⋮ | | | | | |
| 2001-01-01 | Baixo Pantanal | 0 | 23.50009 | 0.08895861 | Pantanais Sul Mato-Grossense |
| 2001-02-01 | Baixo Pantanal | 1 | 23.27970 | 0.25999102 | Pantanais Sul Mato-Grossense |
| 2001-03-01 | Baixo Pantanal | 1 | 22.71204 | 0.26609710 | Pantanais Sul Mato-Grossense |

| meso_name | water_network |
|-----------|---------------|
| Centro Norte De Mato Grosso Do Sul | 86.21000 |
| Pantanais Sul Mato-Grossense | 84.18500 |
| Sudoeste De Mato Grosso Do Sul | 78.20667 |
| Leste De Mato Grosso Do Sul | 79.41250 |

# Hierarchical models

## How to handle variation between groups
## Single-level/non-hierarchical approach

**Varying intercepts**



$$y_i = \boxed{\alpha_{j[i]}} + \beta x_i + \varepsilon_i$$

The average number of cases
observed at the mean temperature in
each region varies,
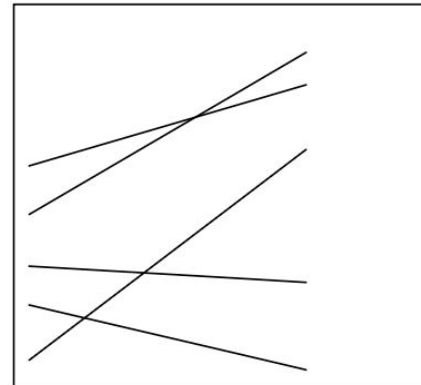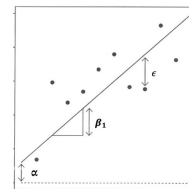but the effect of temperature on
dengue cases is the same

**Varying slopes**



$$y_i = \alpha + \boxed{\beta_{j[i]}} x_i + \varepsilon_i$$

The average number of cases
observed at the mean temperature in
each region is the same, but the
effect of temperature on dengue
cases varies

**Varying intercepts and slopes**



where $j[i]$
is the group
corresponding
to individual $i$

$$y_i = \boxed{\alpha_{j[i]}} + \boxed{\beta_{j[i]}} x_i + \varepsilon_i$$

The average number of cases observed
at the mean temperature in each region
varies, AND the effect of temperature on
dengue cases varies

Gelman and Hill Cambridge University Press 2006

## How to handle variation between groups

(1) **Complete pooling**: assumes there are no differences between the groups. (Equivalent to taking the average number of cases over the entire population or simple linear regression).

$$y_i \sim N\left(\mu_i, \sigma^2\right)$$

$$\mu_i = \alpha + \beta x_i$$

(2) **No pooling**: assumes that each group tells us nothing about any other group. (Equivalent to a separate linear regression for each group or a varying intercept model).

$$y_i \sim N\left(\mu_i, \sigma^2\right)$$

$$\mu_i = \alpha_{j[i]} + \beta x_i$$

### Multilevel/hierarchical approach

(3) **Partial pooling**: pools information across groups by assigning a probability distribution to each group intercept, pulling the group intercept towards the total mean, but allows it to vary by group. (Allows variation of the group-level mean around the total mean).

$$y_i \sim N\left(\mu_i, \sigma^2\right)$$

$$\mu_i = \alpha_{j[i]} + \beta x_i$$

$$\alpha_j \sim N\left(\mu_\alpha, \sigma_\alpha\right)$$

$\mu_\alpha$

$\alpha_j \quad \alpha_j \quad \alpha_j \quad \alpha_j \quad \alpha_j$

Gelman and Hill *Cambridge University Press* 2006

Rowe and Arribas-Bel 2023

## How to handle variation between groups

**No pooling**

**Partial pooling**

Estimated average number of estimated cases in group $j$ ( $\alpha_j$ )

Shrinkage towards the population mean

$\alpha_j$

Number of observations in each group $j$

**Complete pooling**

$\alpha$

Estimates in groups with fewer observations are more variable with higher standard errors.

Estimates in groups with many observations are close to estimate resulting from partial pooling.

Estimates in groups with fewer observations are closer to the complete pooling estimate.

Gelman and Hill Cambridge University Press 2006

## Advantages of hierarchical models

- **Improved estimates for repeated sampling**: When more than one observation arises from the same group (individual, location, or time), then traditional, single-level models either underfit or overfit the data.
- **Improved estimates for imbalance in sampling**: When some groups are sampled more than others, multilevel models prevent over-sampled groups from unfairly dominating inference.
- **Estimates of variation**: multilevel models model variation explicitly, allowing the exploration of individual-level and group-level variation.
- **Avoid averaging, retain variation**: Frequently, scholars pre-average some data to construct variables for a regression analysis. This can be dangerous, because averaging removes variation. Multilevel models allow us to preserve the uncertainty in the original, pre-averaged values, while still using the average to make predictions.

McElreath *Chapman and Hall/CRC* 2020

# Hierarchical models

**Multilevel/hierarchical approach**

Estimates *y*, the model outcome (e.g. case counts), for each observation

$$y_i \sim N\left(\mu_i, \sigma^2\right)$$

$$\mu_i = \alpha_{j[i]} + \beta x_i$$

Estimates $\alpha_j$, the intercept per group (e.g. average number of case counts in each region)

$$\alpha_j \sim N\left(\mu_\alpha, \sigma_\alpha\right)$$

$$\mu_\alpha \sim Normal\,(\,0,\,1.5\,)$$

$$\sigma_\alpha \sim Exponential(\,1\,)$$

Hyperparameters          Hyperpriors

McElreath *Chapman and Hall/CRC* 2020

# Model terms

Monthly Incidence



**Interannual patterns**: is there a common pattern in case incidence every several years, e.g. does every third year usually have more cases on average than July?

**Autocorrelation**: Are years close to each other more likely to have similar values?

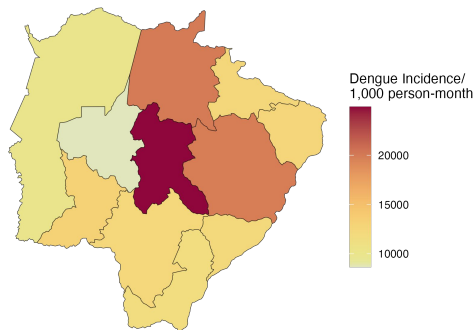$$Y_{s,t} \mid \mu_{s,t}, \theta \sim \text{NegBin}(\mu_{s,t}, \theta)$$

$$\log(\mu_{s,t}) = \alpha + \gamma_{a(t)}$$

$$\boxed{\gamma_a \sim \text{Normal}(0, \tau_a^{-1})}$$

$$\log(\tau_a) = \theta_a$$

$$\theta_a \sim \text{LogGamma}(0.01, 0.01)$$

Here we assume years are *iid*, that is, independent from each other. Other approaches: random walk order 1 or 2.

# Capturing group-level uncertainty - Seasonal



**Seasonality**: is there a common pattern in case incidence every year, e.g. does January usually have more cases on average than July?

**Autocorrelation**: Are years close to each other more likely to have similar values?

$$Y_{s,t} \mid \mu_{s,t}, \theta \sim \mathrm{NegBin}(\mu_{s,t}, \theta)$$

$$\log(\mu_{s,t}) = \alpha + \delta_{m(t)}$$

$$\delta_m - \delta_{m-1} \sim \mathcal{N}(0, \tau^{-1})$$

Here we use a random walk order 1. Other approaches: random walk order 2.

$$\tau \sim \mathrm{Gamma}(0.01, 0.01)$$

# Capturing group-level uncertainty - Spatial

Dengue Incidence/
1,000 person-month

**Regional patterns**: is there a common pattern in case incidence in certain regions, e.g. do the northern regions usually have more cases on average than the southern ones?

**Autocorrelation**: Are regions close to each other more likely to have similar values?



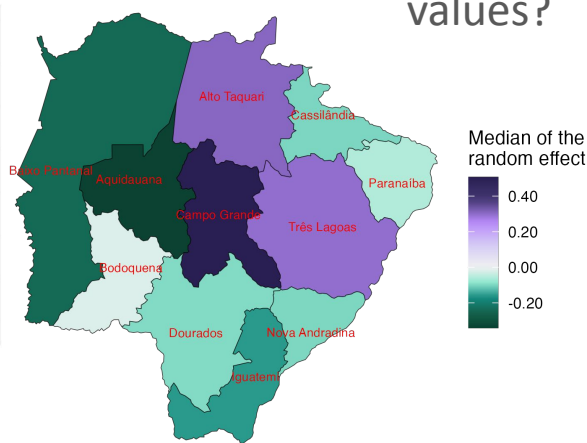Spatial random effects



Spatial random effects

Median of the random effect

FIGURE 8.1: Examples of configurations of areas showing different types of spatial autocorrelation.

Moraga *Chapman and Hall/CRC* 2019

FIGURE 7.7: Left: Areas of the study region. Right: Spatial weight matrix calculated by assuming neighboring areas share a common boundary, and sum of weights for each area.

**ICAR** model

$$u_i | \boldsymbol{u_{-i}} \sim N\left(\bar{u}_{\delta_i}, \frac{\sigma_u^2}{n_{\delta_i}}\right)$$

$$\bar{u}_{\delta_i} = n_{\delta_i}^{-1} \sum_{j \in \delta_i} u_j$$

$\delta_i$ = neighbors of area $i$

$n_{\delta_i}$ = number of neighbors of area $i$

Moraga *Chapman and Hall/CRC* 2019

**g** is an adjacency matrix used to calculate the ICAR (Intrinsic Conditional Auto-Regressive) model used for the prior of the structured spatial effect.
The effect of each area $i$ is normally distributed with a mean equal to the average of its neighbors and a variance decreasing with the number of neighbors.

Here we use a BYM2 prior for the spatial effect. Other approaches include BYM, ICAR, CAR models

$$Y_{s,t} | \mu_{s,t}, \theta \sim \text{NegBin}(\mu_{s,t}, \theta)$$

$$\log(\mu_{s,t}) = \alpha + u_s + v_s$$

Structured spatial effect

Unstructured spatial effect

$$u_s + v_s = \sqrt{\frac{1 - \phi}{\tau}}\, v_s^* + \sqrt{\frac{\phi}{\tau}}\, u_s^*$$

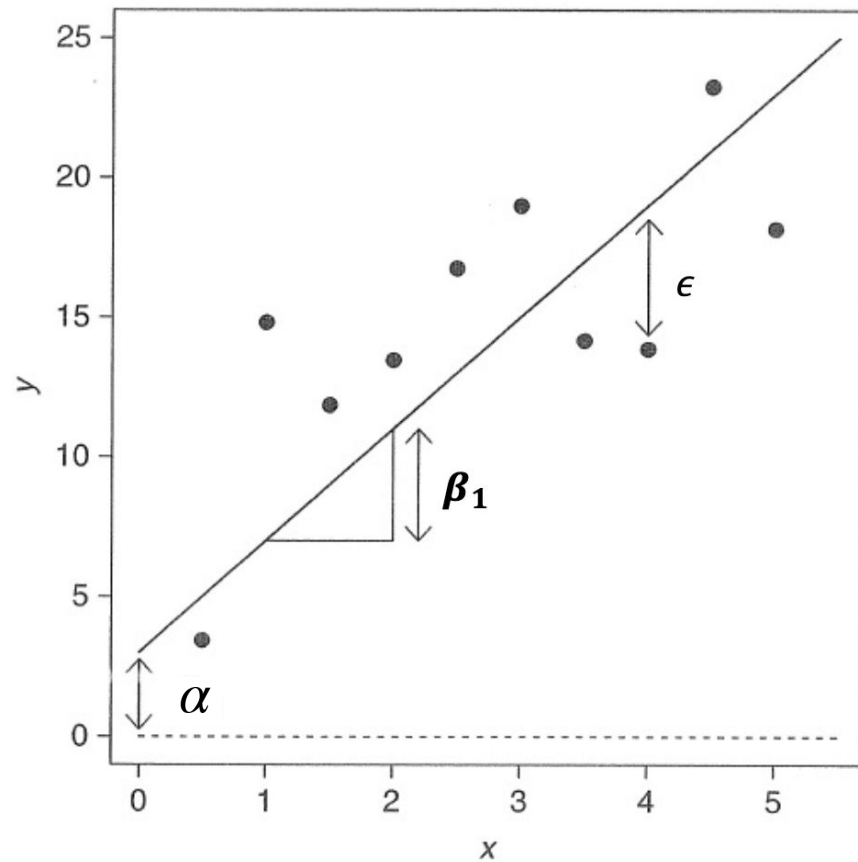$$u_s^* \sim \text{ICAR}(\mathbf{g}) \quad v_s^* \sim \text{Normal}(0, 1)$$

$$\tau \sim \text{PC-Precision}(\sigma = 0.5/0.31, \alpha = 0.01)$$

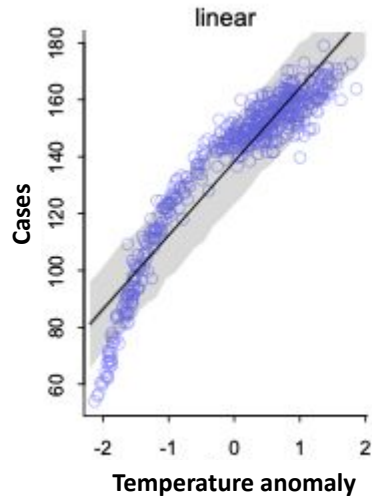$$\phi \sim \text{PC-Mixing}(\phi_0 = 0.5, \alpha = 2/3)$$

$$y_i = \alpha + \beta_1 x_{i1} + \varepsilon_i$$
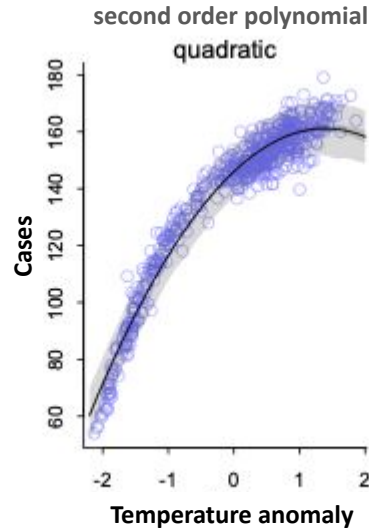
$$y_i \sim N\left(\mu_i, \sigma^2\right)$$

$$\mu_i = \alpha + \beta x_i$$

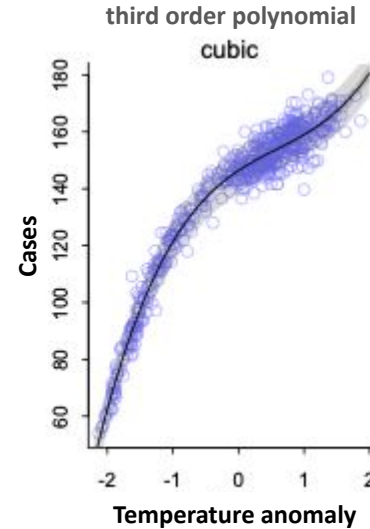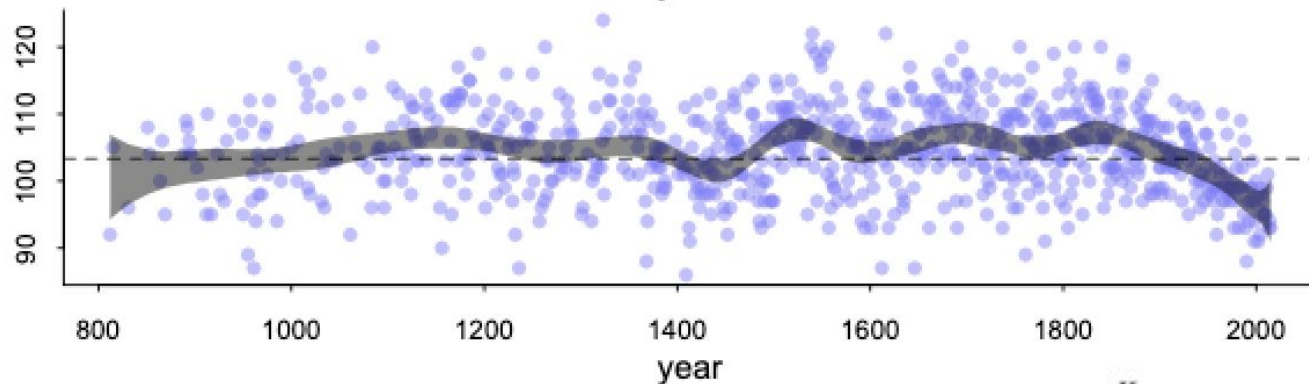# Predictor effects - Non-linear

## Polynomials

$$\mu_i = \alpha + \beta_1 x_{i1}$$

$$\mu_i = \alpha + \beta_1 x_{i1} + \beta_2 x^2_{i1}$$

$$\mu_i = \alpha + \beta_1 x_{i1} + \beta_2 x^2_{i1} + \beta_3 x^3_{i1}$$

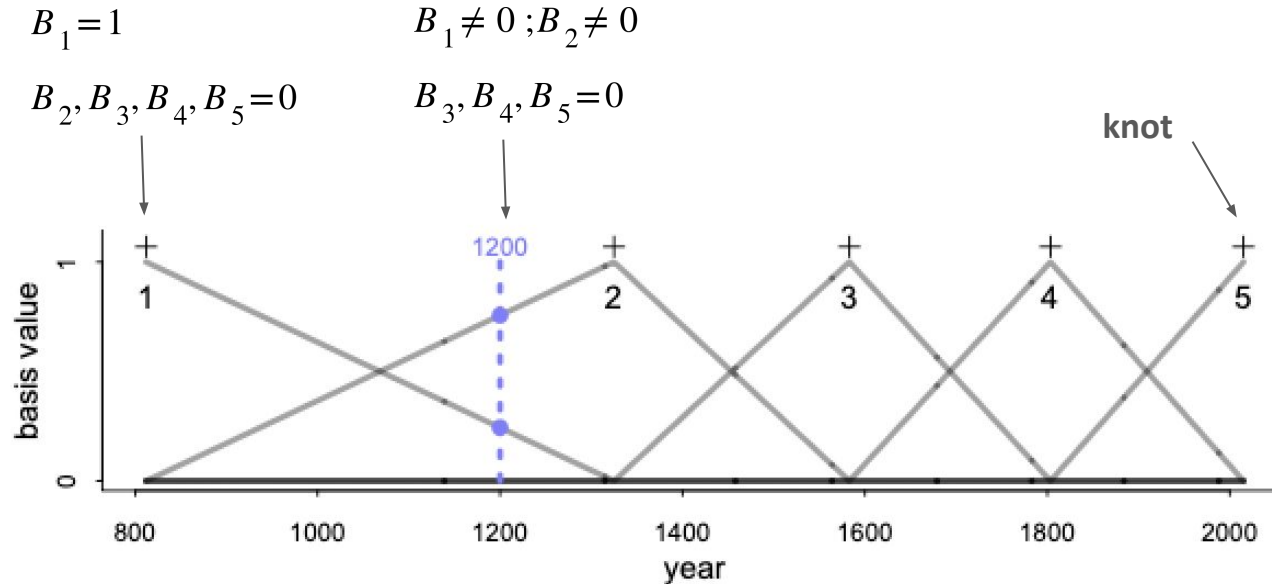McElreath *Chapman and Hall/CRC* 2020

## Splines



Divide the range of *x* variable into parts. Each part has:

**B**: **basis function**: a "synthetic" predictor variable

**w**: **weight parameter**: acts like a slope, adjusting the influence of each basis function on the mean $\mu_i$

$$\mu_i = \alpha + w_1 B_{i,1} + w_2 B_{i,2} + w_3 B_{i,3} + \cdots \longrightarrow \mu_i = \alpha + \sum_{k=1}^{K} w_k B_{k,i}$$

parameter    basis function

McElreath *Chapman and Hall/CRC* 2020

## Splines

$B_1 = 1$

$B_2, B_3, B_4, B_5 = 0$

$B_1 \neq 0 ; B_2 \neq 0$

$B_3, B_4, B_5 = 0$



Divide the range of *x* variable into 4 parts using 5 **knots** placed at even quartiles of the data.

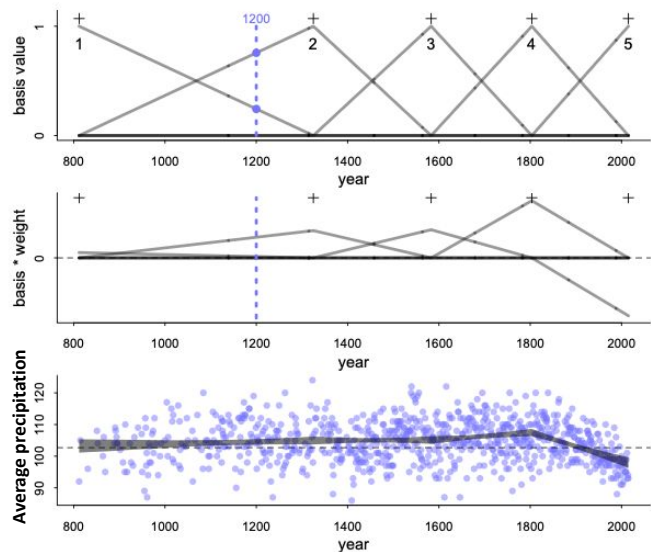*B*: **basis function:** tells you how close you are to each knot.

McElreath *Chapman and Hall/CRC* 2020

# Predictor effects - Non-linear

**Splines**

$$w_1 B_{1200,1} + w_2 B_{1200,2} > 0$$

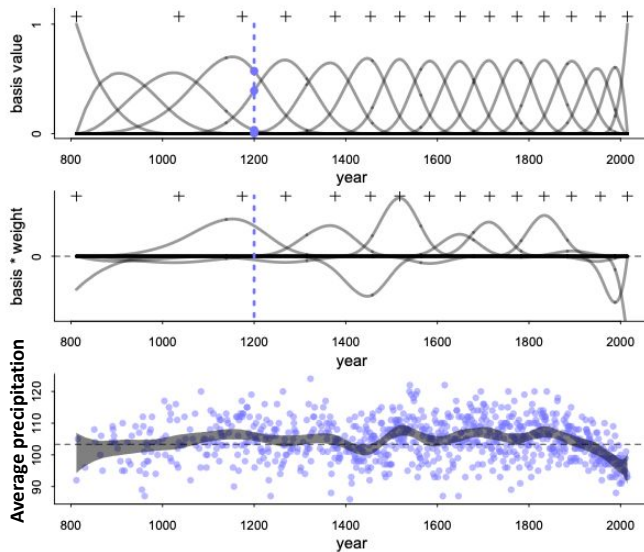*w*: **weight parameters:** are estimated by fitting the model to the data. They can be positive or negative.

Each basis function (*B*) is multiplied by its corresponding weight parameter (*w*).
To predict for a given value of *x*, add the weighted basis functions for that value.

McElreath *Chapman and Hall/CRC* 2020

## How to define spline flexibility?

**Knots = 5; Polynomial degree = 1**
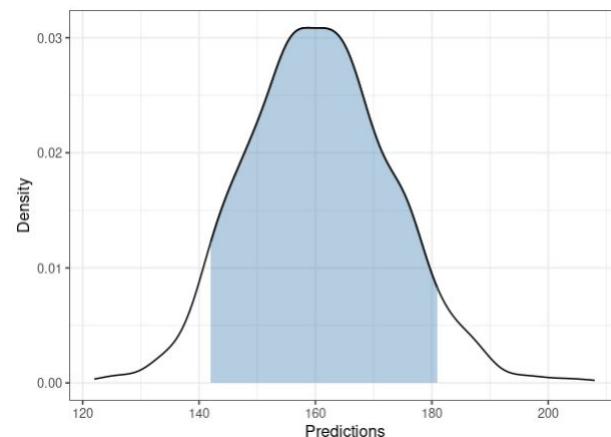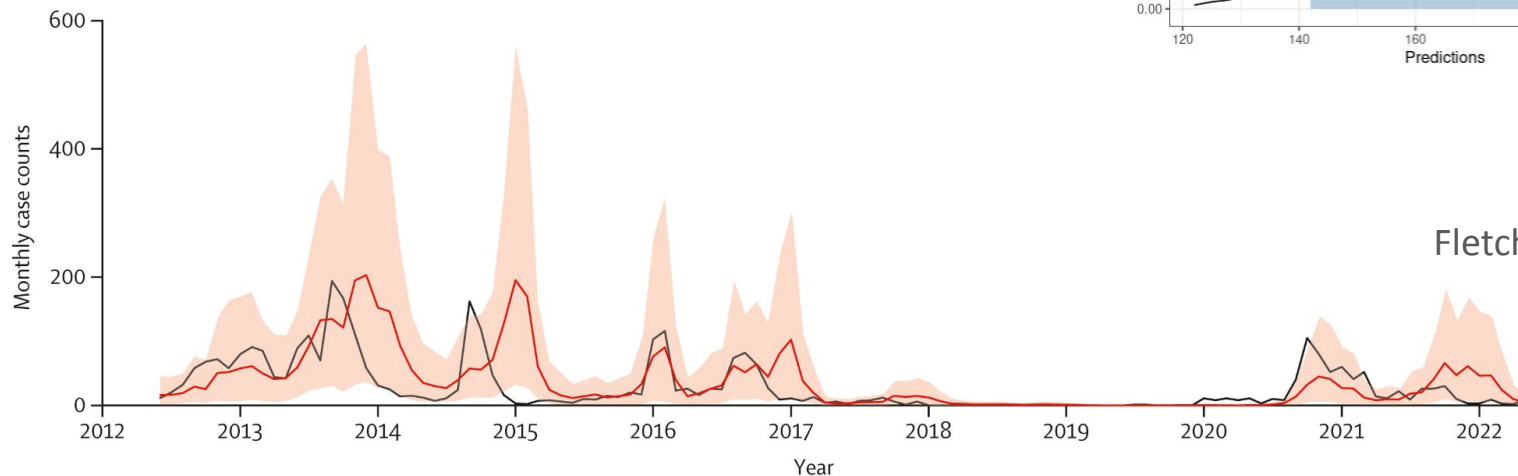
**Knots = 15; Polynomial degree = 3**



- **Number of knots**
- **Placement of knots:** Usually at evenly spaced intervals (equal number of *x* values) or quantiles (equal number of observations)
- **Polynomial degree:** defines how many basis functions combine at each point (value of *x*), that is, how many parameters interact to produce the spline.

McElreath *Chapman and Hall/CRC* 2020

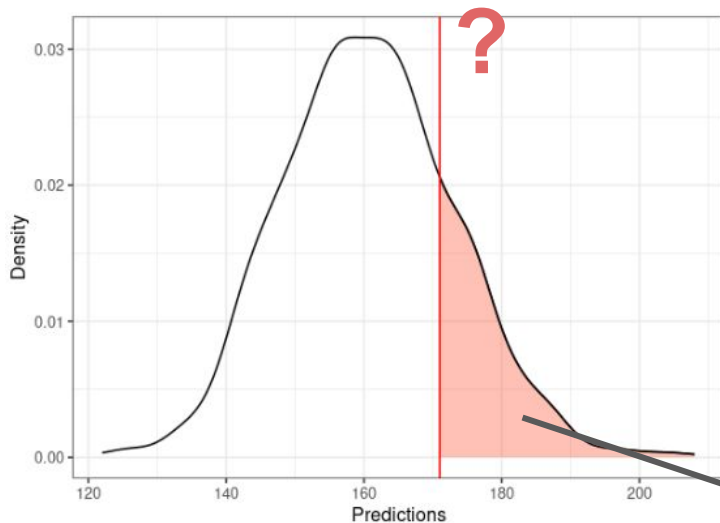# Forecasting for early warning systems

With the fitted model, we can predict case counts (posterior predictive distribution).
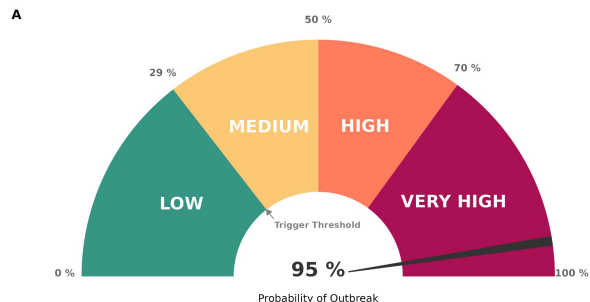
Fletcher et al (2025)

Alternatively, we can communicate our predictions in terms of outbreaks (yes/no) that can trigger the early warning system.
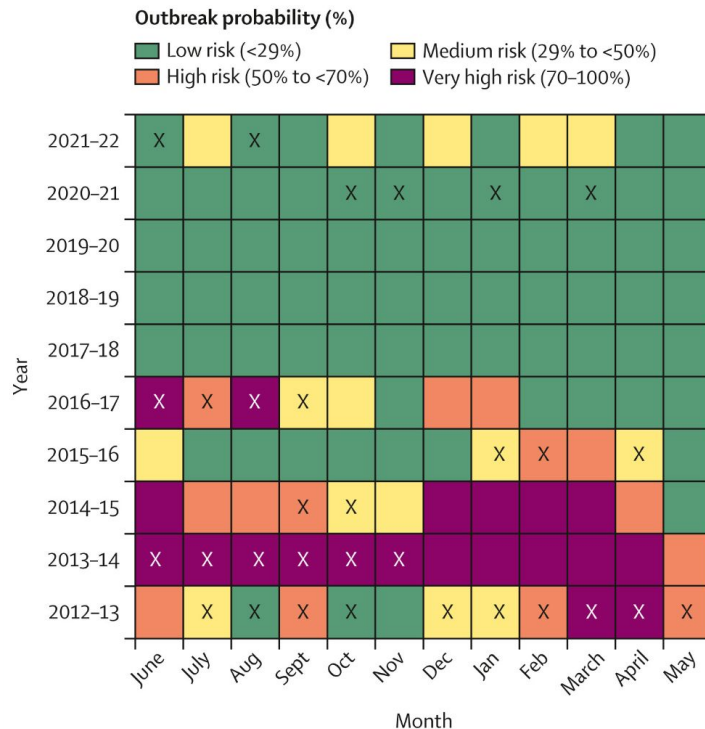


How do we define the <u>outbreak threshold</u>?

- A quantity defined with the stakeholders.

- A certain quantile of the observed cases.

- Mean + φ · SD.

We can calculate the probability of exceeding a threshold: *outbreak probability*

A

We can communicate the probability of outbreaks using different ranges

We can also see how our systems performs using cross-validation

**Source: Fletcher et al. (2025)**

40

# Acknowledgements

The GHR team!!

PROF RACHEL LOWE
LEADING RESEARCHER

DR GIOVENALE MOIRANO
VISITING RESEARCHER

DR MARTIN LOTTO
RESEARCHER

CHLOE FLETCHER
PHD CANDIDATE

# Time for questions

Ania Kawiecki Peralta (ania.kawiecki@bsc.es)

Carles Milà Garcia (carles.milagarcia@bsc.es)

# Useful resources if you want to know more

Gelman, A., & Hill, J. (2006). *Data Analysis Using Regression and Multilevel/Hierarchical Models*.Cambridge University Press.
https://doi.org/10.1017/CBO9780511790942

McElreath, R. (2020). *Statistical Rethinking: A Bayesian Course with Examples in R and STAN* (2nd ed.). Chapman and Hall/CRC.
https://doi.org/10.1201/9780429029608

Moraga, P. (2019). *Geospatial Health Data: Modeling and Visualization with R-INLA and Shiny*. Chapman & Hall/CRC Biostatistics Series.
https://www.paulamoraga.com/book-geospatial/index.html (FREE ONLINE RESOURCE)

Dogucu, M., Ott, M. Q. O., Johnson, A. A. (2021). *Bayes Rules! An Introduction to Applied Bayesian Modeling*.
https://www.bayesrulesbook.com/ (FREE ONLINE RESOURCE)

Rowe, F. and Arribas-Bel, D. (2024) *Spatial Modelling for Data Scientists* https://gdsl-ul.github.io/san/
https://doi.org/10.17605/OSF.IO/8F6XR (FREE ONLINE RESOURCE)

Morris, M. (2019) *Bayesian hierarchical spatial models: Implementing the Besag York Mollié model in stan*. Spatial and Spatio-Temporal
Epidemiology https://doi.org/10.1016/j.sste.2019.100301 (FREE ONLINE RESOURCE here)

Fletcher, C., Moirano, G., Alcayna, T., Rollock, L., Van Meerbeeck, C. J., Mahon, R., Trotman, A., Boodram, L.-L., Browne, T., Best, S.,
Lührsen, D., Diaz, A. R., Dunbar, W., Lippi, C. A., Ryan, S. J., Colón-González, F. J., Stewart-Ibarra, A. M., & Lowe, R. (2025). Compound and
cascading effects of climatic extremes on dengue outbreak risk in the Caribbean: An impact-based modelling framework with long-lag
and short-lag interactions. *The Lancet Planetary Health*, *9*(8), 101279. https://doi.org/10.1016/j.lanplh.2025.06.003 (OPEN SOURCE)