

21th of June 2019

INAR-hidden Markov models to detect and quantify misreported diagnosis in ADHD

Amanda Fernández-Fontelo¹, Alejandra Cabaña², David Moriña², Anna Giménez Palomo³, Pedro Puig²

¹Humboldt-Universität zu Berlin, ²Universitat Autònoma de Barcelona, ³Hospital Clínic de Barcelona

Plan

1 Introduction

- Misreporting
- Motivation: ADHD
- Data
- Goals

2 The model

- Model versions
- Parameter estimation
- Viterbi algorithm: most likely latent (true) sequence
- Model performance

3 Preliminary results

Introduction

Misreporting

Misreporting in data refer to some incident responsible for reporting less (under-reporting) or more (over-reporting) than the actual level of data.

- ▶ Many consequences can derive from misreporting: e.g., inferences emerged from misreported data might be severely biased, drawing an unrealistic picture of the actual problem.
- ▶ Focusing on the public health context, it is well known that some diseases are traditionally under-reported as well as issues such as gender-based violence. However, the over-reporting in some disorders (e.g., ADHD) is extensively documented.

Attention Deficit Hyperactivity Disorder

- ▶ ADHD is a mental disorder usually diagnosed for the first time in school-aged children. Although this is one of the most common mental disorders in children, it remains poorly understood.
- ▶ Some of the more frequently documented reasons related to this systematic misdiagnosis in ADHD are:
 - ▶ relative age of school-age children.
 - ▶ symptoms are drastically different between boys and girls.
 - ▶ disagreement between diagnosis protocols.
 - ▶ etc.

Ford-Jones, P.C. (2015). Misdiagnosis of attention deficit hyperactivity disorder: “Normal behaviour” and relative maturity, *Paediatric Child Health*, 20(4): 200–202.

Data

Data are the monthly number of visits to CSMA from 2013 to 2018, and come from “Agència de Qualitat i Avaluació Sanitàries de Catalunya (AQuAS)”¹:

- ▶ Geographic area (comarca and municipi of Catalonia).
 - Barcelonès, Vallès Occidental and Baix Llobregat.
- ▶ Age range (< 1 year, then 5-year intervals from 1 year to 79).
 - Re-categorized in children (< 14), adolescents (15 – 24) and adults (> 24).
- ▶ Gender.
- ▶ Relevance of ADHD in the visit: is this disease the first reason for visiting?
 - Considered cases where ADHD is the primary reason for visiting.

¹<http://aquas.gencat.cat/ca/inici>

Goals

The goal of this work consists of providing a novel tool able to detect and quantify the misreporting in both directions (under- and over-reporting):

- ▶ Apply the tool to detect misreporting in ADHD cases in Catalonia according to the geographical area, gender, and age to provide clinicians a more objective measure to such misdiagnosis problem.

The model



Model I: considering only under-reporting

Consider a **latent process** X_n with the following Poisson(λ)-INAR(1) structure:

$$X_n = \alpha \circ X_{n-1} + W_n(\lambda),$$

where $\alpha \in (0, 1)$. $E(X_n) = V(X_n) = \lambda/(1 - \alpha) = \mu_X$. The operator \circ is the binomial thinning such that: $\alpha \circ X_{n-1} = \sum_{i=1}^{X_{n-1}} Z_i$, where Z_i are i.i.d Bernoulli(α).

Consider the following **observed and potentially under-reported process** Y_n : X_n with probability $1 - \omega$, or $q \circ X_n$ with probability ω .

Fernández-Fontelo, A., Cabaña, A., Puig, P. and Moríña, D. (2016). Under-reported data analysis with INAR-hidden Markov chains. *Statistics in Medicine*, 35(26): 4875-4890.

Model II: considering under-reporting under a more sophisticated correlation structure

Imagine now that the under-reporting indicators $\{\mathbf{1}_n, n \geq 1\}$ are also serially dependent.

- ▶ Take the simplest scheme: a binary discrete-time Markov chain.

Now, the **observed and potentially under-reported process** $\{Z_n\}$ is: X_n with probability $1 - \omega$, or $q \circ X_n$ with probability ω , where $P(\mathbf{1}_{n+k} = 1_{n+k} | \mathbf{1}_n = 1_n) \neq P(\mathbf{1}_{n+k} = 1_{n+k})$.

- ▶ One additional parameter is included in the model to accommodate the correlation among the under-reporting states.

Model III: Considering both under-reporting and over-reporting

Still consider that the latent process X_n is a $\text{Poisson}(\lambda)$ -INAR(1) model, but now the processes Y_n (and Z_n) can be misreported (under- or over-reported):

$$Y_n = \begin{cases} X_n & 1 - \omega \\ \theta \diamond X_n & \omega, \end{cases}$$

where $\theta = (\phi_1, \phi_2)$ and $\theta \diamond X_n$ is called the fattening-thinning operator:
 $\theta \diamond X_n | X_n = x_n = \sum_{j=1}^{x_n} W_j$:

$$P(W = j) = \begin{cases} 0 & 1 - \phi_1 - \phi_2 \\ 1 & \phi_1 \\ 2 & \phi_2 \end{cases}.$$

Under-reporting or over-reporting ?

- ▶ To distinguish between under-reporting, no misreporting or over-reporting, the following can easily be computed once the model is estimated:

under-reporting	$\phi_1 + 2\phi_2 < 1$
no misreporting	$\phi_1 + 2\phi_2 = 1$
over-reporting	$\phi_1 + 2\phi_2 > 1$

- ▶ Notice that when $\phi_2 = 0$, the model results in the versions I and II, which only accounts for under-reporting or no misreporting.

Parameter estimation: moment-based method

The marginal distribution of the observed process (Y_n or Z_n) is essential to compute the moment-based estimates of the model:

$$Y_n = \begin{cases} \text{Poisson}(\mu_X) & 1 - \omega, \\ \text{2nd-order Hermite}(\mu_X\phi_2, \mu_X(1 - \phi_1 - \phi_2)) & \omega. \end{cases}$$

- 1 Fit the mixture above to obtain estimates of $\hat{\omega}$, $\hat{\mu}_X$, $\hat{\phi}_1$ and $\hat{\phi}_2$.
- 2 Use the theoretical expression of the ACF (ρ_Y and ρ_Z) to estimate α .
- 3 Using $\hat{\mu}_X$ (step 1) and $\hat{\alpha}$ (step 2), λ can be easily estimated.

If G_1 and G_2 are independent Poisson distributions with parameters a and b , $H = G_1 + 2G_2 \sim \text{2nd-order Hermite}(a, b)$.

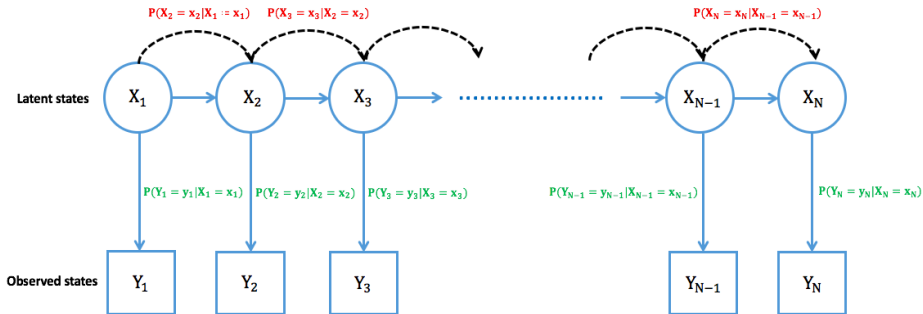
Parameter estimation: likelihood-based method

Since the likelihood functions of the processes Y_n and Z_n are directly intractable (HMC with an infinite number of states), the forward algorithm is a reasonable choice to compute such functions.

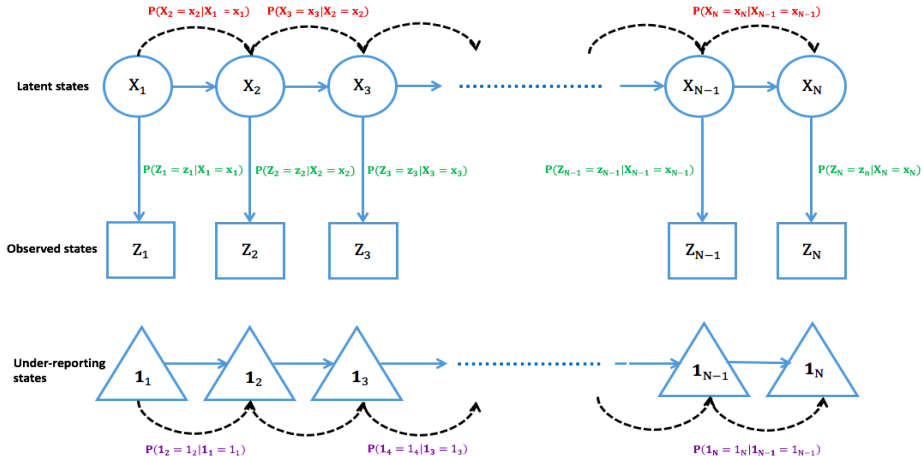
The likelihood is then achieved by $P(Y_{1:N} = y_{1:N}) = \sum_{x_N = \frac{y_N}{2}, \mathbf{1}_N}^{\infty} \gamma_N(y_{1:N}, x_N, \mathbf{1}_n)$, where:

$$\begin{aligned} \gamma_n(y_{1:n}, x_n, \mathbf{1}_n) &= \overbrace{P(Y_n = y_n | X_n = x_n, \mathbf{1}_n)}^{\text{emission probabilities}} \sum_{x_{n-1} = \frac{y_{n-1}}{2}, \mathbf{1}_{n-1}}^{\infty} \underbrace{P(X_n = x_n | X_{n-1} = x_{n-1})}_{\text{transition probabilities}} \\ &\times \underbrace{P(\mathbf{1}_n = \mathbf{1}_n | \mathbf{1}_{n-1} = \mathbf{1}_{n-1})}_{\text{transition probability matrix}} \gamma_{n-1}(y_{1:n-1}, x_{n-1}, \mathbf{1}_{n-1}). \end{aligned}$$

Visual model



Visual model



Forward probabilities

- While the transition probabilities are easily computed through the Poisson(λ)-INAR(1) cpmf, the emission probabilities are trickier:

model	$P(Y_n = y_n X_n = x_n, \mathbf{1}_n = \mathbf{1}_n)$
Y_n	$\begin{cases} 0 & \text{if } x_n < \frac{y_n}{2} \\ (1 - \omega) + \omega p_n & \text{if } y_n = x_n \\ \omega p_n & \text{if } x_n \geq \frac{y_n}{2} \end{cases}$
Z_n	$\begin{cases} 0 & \text{if } x_n < \frac{y_n}{2} \\ 0 & \text{if } \mathbf{1}_n = 0, x_n < \frac{y_n}{2} \\ p_n & \text{if } \mathbf{1}_n = 1, x_n \geq \frac{y_n}{2} \\ 1 & \text{if } \mathbf{1}_n = 0, x_n = y_n \end{cases}$

- p_n can be computed through the following recursive relation:

$$p_n = \frac{1}{n(1 - \phi_1 - \phi_2)} [\phi_1(x_n - (n - 1))p_{n-1} + \phi_2(x_n - (n - 2))p_{n-2}]$$

Viterbi algorithm: most likely latent (true) sequence

The Viterbi algorithm ² is used to know the latent chain that maximises $P(X_{1:n}|Y_{1:n}) = \frac{P(X_{1:n}, Y_{1:n})}{P(Y_{1:n})}$ (assuming all parameters are known):

- ▶ Since $P(Y_{1:n})$ does not depend on X_n , it is enough to maximise the probability $P(X_{1:n}, Y_{1:n})$.
- ▶ The most likely chain of latent states is obtained as:

$$X^* = \arg \max_X P(X_{1:n}, Y_{1:n}).$$

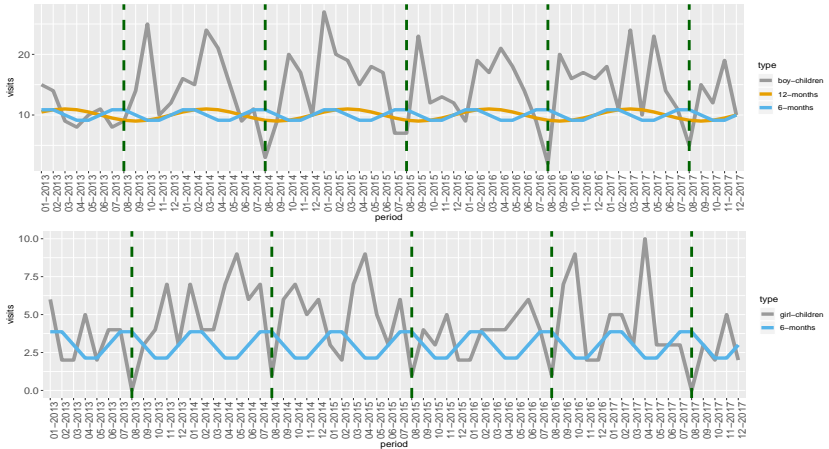
²Viterbi AJ. Error bounds for convolutional codes and an asymptotically optimum decoding algorithm. *IEEE Transactions on Information Theory* Apr 1967, 13, 260-269.

Preliminary results

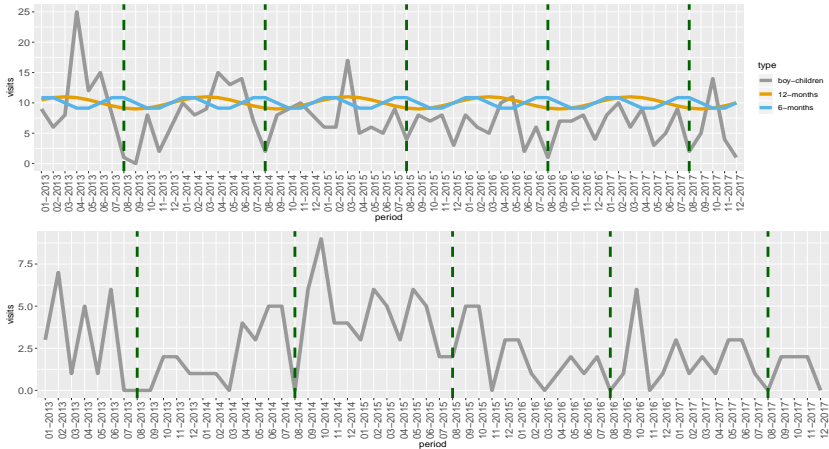
Exploratory analysis

- ▶ Nearly 60 series among Barcelonès, Vallès Occidental and Baix Llobregat. Some of them are rejected because of low counts and high frequency of zeros (convergence problems in MLE).
- ▶ Many of them (especially boy-children) show clear annual and semiannual seasonal patterns (e.g., vacation periods).
- ▶ Example of boy-children and girl-children in Hospitalet de Llobregat (Barcelonès) and Terrassa (Vallès Occidental).

Seasonal patterns in Hospitalet de Llobregat



Seasonal patterns in Terrassa



Proposed models

area	gender	model	misreporting
Hospitalet de Llob.	boy	$X_n \sim \text{Poisson} \left(e^{2.3224 + 0.1742 \sin \frac{2\pi n}{12} + 0.1314 \cos \frac{2\pi n}{12}} \right)$ $Y_n = \begin{cases} X_n & 0.1325 \\ \hat{\theta} = (0.3122, 0.5632) \diamond X_n & 0.8675 \end{cases}$	$0.3122 + 20.5632 > 1$
	girl	$Y_n \sim \text{Poisson} \left(e^{1.4272 - 0.2812 \sin \frac{2\pi n}{6} + 0.0118 \cos \frac{2\pi n}{6}} \right)$	-
Terrassa	boy	$X_n \sim \text{Poisson} \left(e^{2.8208 + 0.2687 \sin \frac{2\pi n}{12} + 0.2040 \cos \frac{2\pi n}{12}} \right)$ $Y_n = \begin{cases} X_n & 0.0275 \\ \hat{\theta} = (0.0629, 0.1776) \diamond X_n & 0.9725 \end{cases}$	$0.0629 + 20.1776 < 1$
	girl	$Y_n \sim \text{Poisson}(4.2500)$	-

To Do

- ▶ Series should be conveniently validated through pseudo-residuals, and the most likely sequence of latent state should be provided (with bands).
- ▶ The reminder series of Barcelonès, Vallès Occidental, and Baix Llobregat should be analyzed appropriately.
- ▶ Seasonal patterns have to be validated with clinicians. Furthermore, further discussions with psychiatrists are needed to confirm the differences in misreporting depending on geographical areas.

[illegible]