

Zero-inflated negative binomial mixed models

for predicting number of wildfires

M. Bugallo¹, M.D. Esteban¹, D. Morales¹, M.F. Marey-Pérez²

¹Center of Operations Research, Miguel Hernández University of Elche, Spain

²Higher Polytechnic School of Engineering, Santiago de Compostela University, Spain

Stochastic and Physical Monitoring Systems (SPMS 2024)
Group of Applied Mathematics and Stochastics (GAMS)
Department of Mathematics, Czech Technical University in Prague
Sokol Dobručovice, Czech Republic, June 20-24, 2024



UNIVERSITAS
Miguel Hernández
RESEARCH INSTITUTE



UNIVERSITAS
Miguel Hernández

Overview

Data and problem of interest

Model and predictors

Bootstrap-based inference

Model-based statistical analysis

Wildfire forecasting and error measurements

Conclusions

Appendix: ML-Laplace approximation algorithm

Data and problem of interest

Model and predictors

Bootstrap-based inference

Model-based statistical analysis

Wildfire forecasting and error measurements

Conclusions

Appendix: ML-Laplace approximation algorithm

Wildfires

- **Figure 1** shows the distribution and extent of burnt areas in Europe and the Mediterranean in 2021.
- Red circles represent fire events.

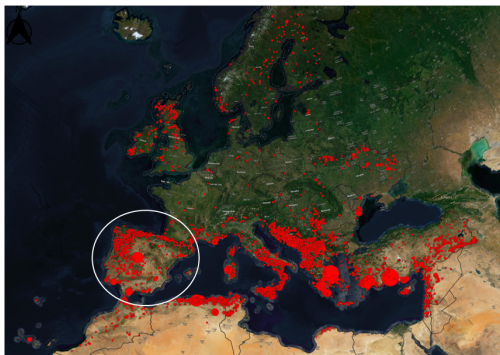


Figure 1: Wildfires in 2021.

Target variable

- We collect data on Spanish forest fires, at the provincial level for all months 2002-2015, with a total of 216,538 events.
- The dependent variable of interests, y_{ijk} , counts the **number of wildfires** in year i , month j and province k .
- We have $D = IJK = 8400$ domains, defined by the crossings of years ($I = 14$), months ($J = 12$) and provinces ($K = 50$).
- There are 951 domains where no forest fire was recorded.

2002	2003	2004	2005	2006	2007	2008
107	112	63	51	90	52	60
2009	2010	2011	2012	2013	2014	2015
47	89	44	55	60	77	44

Table 1: Number of zeros per year.

Target variable

Jan.	Feb.	March	April	May	June
177	112	52	40	22	20
July	Aug.	Sept.	Oct.	Nov.	Dec.
5	6	9	59	206	243

Table 2: Number of zeros per month.

Jan.	Feb.	March	April	May	June
6372	19270	30432	16701	11844	18494
July	Aug.	Sept.	Oct.	Nov.	Dec.
28912	37742	26124	12033	4343	4271

Table 3: Number of wildfires per month.

Auxiliary variables

- *year3* is an explanatory variable, with three categories indicating groups of years.
- *year3.1* corresponds to the years 2002-2006. It contains 3000 domains and 423 zeros (14.1%).
 - The average ($50 \times 3 \times 5$) of total burned areas in July-September was 668.51 Ha, with a total of 48,518 active fires.
- *year3.2* is for the years 2007-2012. It contains 3600 domains and 347 zeros were observed (9.64%).
 - The average of total burned areas in July-September was 333.56 Ha, with a total of 30,020 active fires.
- *year3.3* is for years 2013-2015. It contains 1800 domains, such that 181 zeros were observed (10.06%)
 - The average of total burned areas in July-September was 265.48 Ha, with a total of 14,240 active fires.

Auxiliary variables

Variable	Description	Units
<i>e</i>	average vapor pressure	tenths of hPa
<i>hr</i>	average relative humidity	tenths of mm
<i>n.fog</i>	foggy days	% days
<i>n.gra</i>	hail days	% days
<i>n.llu</i>	rainy days	% days
<i>n.nie</i>	snowy days	% days
<i>np.001</i>	precipitation ≥ 0.1 mm	% days
<i>np.010</i>	precipitation ≥ 1 mm	% days
<i>np.300</i>	precipitation ≥ 30 mm	% days
<i>nt.00</i>	min. temperature $\leq 0^{\circ}\text{C}$	% days
<i>nt.30</i>	max. temperature $\geq 30^{\circ}\text{C}$	% days
<i>n.tor</i>	storm days	% days
<i>nw.55</i>	wind speed ≥ 55 km/h	% days
<i>nw.91</i>	wind speed ≥ 91 km/h	% days
<i>p.max</i>	max, daily rainfall	mm
<i>p.mes</i>	total precipitation	mm

Table: Description of the auxiliary variables.

Auxiliary variables

Variable	Description	Units
<i>q.mar</i>	mean sea-level pressure	KPa
<i>q.max</i>	max. absolute pressure	KPa
<i>q.med</i>	average pressure	KPa
<i>q.min</i>	max. min. pressure	KPa
<i>ta.max</i>	absolute max. temperature	°C
<i>ta.min</i>	absolute min. temperature	°C
<i>ti.max</i>	lowest max. temperature	°C
<i>tm.max</i>	average max. temperature	°C
<i>tm.mes</i>	average temperature	°C
<i>tm.min</i>	average min. temperature	°C
<i>ts.min</i>	highest min. temperature	°C
<i>w.med</i>	average speed elaborated from 07, 13, 18 UTC	km/h
<i>unemp</i>	unemployment rate	%
<i>year3</i>	year group variable	-

Table: Description of the auxiliary variables.

Data and problem of interest

Model and predictors

Bootstrap-based inference

Model-based statistical analysis

Wildfire forecasting and error measurements

Conclusions

Appendix: ML-Laplace approximation algorithm

Variables

Target variables

- y_{ijk} is a count variable taking values on $\mathbb{N} \cup \{0\}$,
 $i \in \mathbb{I} = \{1, \dots, I\}$, $j \in \mathbb{J} = \{1, \dots, J\}$, $k \in \mathbb{K} = \{1, \dots, K\}$.
- y_{ijk} is the number of wildfires in year i , month j and province k .
- z_{ijk} is a zero-inflation latent (non observable) variable.
- $D = IJK$ is the total number of domains.

Explanatory variables

- $\mathbf{x}_{1,ijk} = (x_{1,ijk1}, \dots, x_{1,ijkq_1})$ and $\mathbf{x}_{2,ijk} = (x_{2,ijk1}, \dots, x_{2,ijkq_2})$
 are $1 \times q_1$ and $1 \times q_2$ row vectors with explanatory variables.

$$\mathbf{y}_{jk} = \text{col}_{1 \leq i \leq I} (y_{ijk}), \mathbf{z}_{jk} = \text{col}_{1 \leq i \leq I} (z_{ijk}), \mathbf{y} = \text{col}_{1 \leq j \leq J} \left(\text{col}_{1 \leq k \leq K} (\mathbf{y}_{jk}) \right), \mathbf{z} = \text{col}_{1 \leq j \leq J} \left(\text{col}_{1 \leq k \leq K} (\mathbf{z}_{jk}) \right)$$

$$\mathbf{X}_{a,jk} = \text{col}_{1 \leq k \leq K} (\mathbf{x}_{a,ijk}), \mathbf{X}_a = \text{col}_{1 \leq j \leq J} \left(\text{col}_{1 \leq k \leq K} (\mathbf{X}_{a,jk}) \right), a = 1, 2.$$

Random effects

- $u_{1,j}, u_{1,k}, u_{2,j}, u_{2,k}, j \in \mathbb{J}, k \in \mathbb{K}$, independent and $N(0, 1)$.
- Define the vector $\mathbf{u} = (\mathbf{u}'_1, \mathbf{u}'_2)'$, where

$$\begin{aligned} \mathbf{u}_{1,jk} &= (u_{1,j}, u_{1,k})', \quad \mathbf{u}_{2,jk} = (u_{2,j}, u_{2,k})', \quad \mathbf{u}_{jk} = (\mathbf{u}'_{1,jk}, \mathbf{u}'_{2,jk})', \\ \mathbf{u}_1 &= \operatorname{col}_{1 \leq j \leq J} \left(\operatorname{col}_{1 \leq k \leq K} (\mathbf{u}_{1,jk}) \right) \sim N_{2JK}(0, \mathbf{I}), \\ \mathbf{u}_2 &= \operatorname{col}_{1 \leq j \leq J} \left(\operatorname{col}_{1 \leq k \leq K} (\mathbf{u}_{2,jk}) \right) \sim N_{2JK}(0, \mathbf{I}), \quad \mathbf{u} = (\mathbf{u}'_1, \mathbf{u}'_2)'. \end{aligned}$$

The AZINB11 model (NB2 parameterization)

- $(z_{ijk}, y_{ijk}) \sim \text{AZINB11}, i \in \mathbb{I}, j \in \mathbb{J}, k \in \mathbb{K}$, if

$$z_{ijk} \sim \text{Bern}(p_{ijk}), \quad P(y_{ijk} = 0 / z_{ijk} = 1) = 1, \quad y_{ijk} |_{z_{ijk}=0} \sim \text{NB}(r, \mu_{ijk}), \text{ i.e.}$$

$$P(y_{ijk} = t / z_{ijk} = 0) = \frac{\Gamma(t+r)}{\Gamma(t+1)\Gamma(r)} \left(\frac{\mu_{ijk}}{r + \mu_{ijk}} \right)^t \left(\frac{r}{r + \mu_{ijk}} \right)^r, \quad t \in \mathbb{N} \cup \{0\},$$

where $p_{ijk} \in (0, 1), r > 0, \mu_{ijk} > 0$.

The model

To complete the definition of the AZINB11 model, we assume

- The link functions are

$$\text{logit}(p_{ijk}) = \log \frac{p_{ijk}}{1 - p_{ijk}} = \mathbf{x}_{1,ijk} \boldsymbol{\beta}_1 + \phi_{11} u_{1,j} + \phi_{12} u_{1,k},$$

$$\log(\mu_{ijk}) = \mathbf{x}_{2,ijk} \boldsymbol{\beta}_2 + \phi_{21} u_{2,j} + \phi_{22} u_{2,k}.$$

- Conversely, for $i \in \mathbb{I}$, $j \in \mathbb{J}$, $k \in \mathbb{K}$, we have

$$p_{ijk} = \frac{\exp\{\mathbf{x}_{1,ijk} \boldsymbol{\beta}_1 + \phi_{11} u_{1,j} + \phi_{12} u_{1,k}\}}{1 + \exp\{\mathbf{x}_{1,ijk} \boldsymbol{\beta}_1 + \phi_{11} u_{1,j} + \phi_{12} u_{1,k}\}},$$

$$\mu_{ijk} = \exp\{\mathbf{x}_{2,ijk} \boldsymbol{\beta}_2 + \phi_{21} u_{2,j} + \phi_{22} u_{2,k}\}.$$

- The vectors $(y_{ijk}, z_{ijk})'$, $i \in \mathbb{I}$, $j \in \mathbb{J}$, $k \in \mathbb{K}$, are independent conditioned to \mathbf{u} .

The likelihood

Model parameters

- $\boldsymbol{\theta}_1 = (\boldsymbol{\beta}'_1, \phi_{11}, \phi_{12})'$, $\boldsymbol{\theta}_2 = (\boldsymbol{\beta}'_2, \phi_{21}, \phi_{22})'$, $\boldsymbol{\theta} = (\boldsymbol{\theta}'_1, \boldsymbol{\theta}'_2)'$,

Conditioned likelihood factors

$$\begin{aligned}
 P(y_{ijk} | \mathbf{u}_{jk}; \boldsymbol{\theta}) &= (1 + \exp\{\mathbf{x}_{1,ijk}\boldsymbol{\beta}_1 + \phi_{11}u_{1,j} + \phi_{12}u_{1,k}\})^{-1} \\
 &\cdot \left\{ \xi_{ijk} \left[\exp\{\mathbf{x}_{1,ijk}\boldsymbol{\beta}_1 + \phi_{11}u_{1,j} + \phi_{12}u_{1,k}\} \right. \right. \\
 &+ \left. \left. \exp\left\{r \log r - r \log(r + \exp\{\mathbf{x}_{2,ijk}\boldsymbol{\beta}_2 + \phi_{21}u_{2,j} + \phi_{22}u_{2,k}\})\right\} \right] \right\} \\
 &+ (1 - \xi_{ijk}) \exp\left\{y_{ijk}(\mathbf{x}_{2,ijk}\boldsymbol{\beta}_2 + \phi_{21}u_{2,j} + \phi_{22}u_{2,k})\right. \\
 &- (y_{ijk} + r) \log(r + \exp\{\mathbf{x}_{2,ijk}\boldsymbol{\beta}_2 + \phi_{21}u_{2,j} + \phi_{22}u_{2,k}\}) \\
 &\left. + r \log r + \log \frac{\Gamma(y_{ijk} + r)}{\Gamma(y_{ijk} + 1)\Gamma(r)} \right\}.
 \end{aligned}$$

The likelihood

Conditioned likelihood

$$P(\mathbf{y}|\mathbf{u}; \boldsymbol{\theta}) = \prod_{j=1}^J \prod_{k=1}^K P(\mathbf{y}_{jk}|\mathbf{u}_{jk}; \boldsymbol{\theta}), \quad P(\mathbf{y}_{jk}|\mathbf{u}_{jk}; \boldsymbol{\theta}) = \prod_{i=1}^I P(y_{ijk}|\mathbf{u}_{jk}; \boldsymbol{\theta}).$$

Marginal likelihood

$$P(\mathbf{y}; \boldsymbol{\theta}) = \prod_{j=1}^J \prod_{k=1}^K \int_{\mathbb{R}^4} \prod_{i=1}^I P(y_{ijk}|\mathbf{u}_{jk}; \boldsymbol{\theta}) f_{N_4(0, I)}(\mathbf{u}_{jk}) d\mathbf{u}_{jk}.$$

Log-likelihood function

$$\ell(\boldsymbol{\theta}; \mathbf{y}) = \sum_{j=1}^J \sum_{k=1}^K \log \int_{\mathbb{R}^4} \prod_{i=1}^I P(y_{ijk}|\mathbf{u}_{jk}; \boldsymbol{\theta}) f_{N_4(0, I)}(\mathbf{u}_{jk}) d\mathbf{u}_{jk}.$$

ML estimators

- Maximum likelihood (ML) estimators

$$\hat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \ell(\boldsymbol{\theta}; \mathbf{y}), \quad \Theta = \mathbb{R}^{q_1+q_2} \times \mathbb{R}_+^4, \quad \mathbb{R}_+ = (0, \infty).$$

- As $\ell(\boldsymbol{\theta}; \mathbf{y})$ contains integrals in \mathbb{R}^4 , we apply a Laplace approximation algorithm to
 - maximize $\ell(\boldsymbol{\theta}; \mathbf{y})$
 - calculate the ML estimators of the model parameters: $\hat{\boldsymbol{\theta}}$,
 - obtain mode predictors of random effects: $\hat{\mathbf{u}}$.

Predictors

Prediction

- We are interested in predicting expected counts

$$\begin{aligned}\mu_{yijk} &= E[y_{ijk} | \mathbf{u}_{jk}] = (1 - p_{ijk}(\boldsymbol{\theta}_1, \mathbf{u}_{1,jk})) \mu_{ijk}(\boldsymbol{\theta}_2, \mathbf{u}_{2,jk}) \\ p_{ijk}(\boldsymbol{\theta}_1, \mathbf{u}_{1,jk}) &= \frac{\exp\{\mathbf{x}_{1,ijk} \boldsymbol{\beta}_1 + \phi_{11} u_{1,j} + \phi_{12} u_{1,k}\}}{1 + \exp\{\mathbf{x}_{1,ijk} \boldsymbol{\beta}_1 + \phi_{11} u_{1,j} + \phi_{12} u_{1,k}\}}, \\ \mu_{ijk}(\boldsymbol{\theta}_2, \mathbf{u}_{2,jk}) &= \exp\{\mathbf{x}_{2,ijk} \boldsymbol{\beta}_2 + \phi_{21} u_{2,j} + \phi_{22} u_{2,k}\}.\end{aligned}$$

- The plug-in predictor of μ_{yijk} is

$$\hat{\mu}_{yijk}^{in} = (1 - p_{ijk}(\hat{\boldsymbol{\theta}}_1, \hat{\mathbf{u}}_{1,jk})) \mu_{ijk}(\hat{\boldsymbol{\theta}}_2, \hat{\mathbf{u}}_{2,jk})$$

- We are **not** interested in predicting expected NB counts μ_{ijk} by using the EBP $E[\mu_{ijk} | \mathbf{y}]$.

Data and problem of interest

Model and predictors

Bootstrap-based inference

Model-based statistical analysis

Wildfire forecasting and error measurements

Conclusions

Appendix: ML-Laplace approximation algorithm

- 1 Calculate the ML estimate $\hat{\theta}$. Run **Algorithm B1**.
- 2 Repeat B times ($b = 1, \dots, B$):
 - (a) Generate $u_{1,j}^{*(b)} \sim N(0, 1)$, $u_{1,k}^{*(b)} \sim N(0, 1)$, $u_{2,j}^{*(b)} \sim N(0, 1)$ and $u_{2,k}^{*(b)} \sim N(0, 1)$. Calculate

$$p_{ijk}^{*(b)} = \frac{\exp \{ \mathbf{x}_{1,ijk} \hat{\beta}_1 + \hat{\phi}_{11} u_{1,j}^{*(b)} + \hat{\phi}_{12} u_{1,k}^{*(b)} \}}{(1 + \exp \{ \mathbf{x}_{1,ijk} \hat{\beta}_1 + \hat{\phi}_{11} u_{1,j}^{*(b)} + \hat{\phi}_{12} u_{1,k}^{*(b)} \})},$$

$$\mu_{ijk}^{*(b)} = \exp \{ \mathbf{x}_{2,ijk} \hat{\beta}_2 + \hat{\phi}_{21} u_{2,j}^{*(b)} + \hat{\phi}_{22} u_{2,k}^{*(b)} \}.$$

- (b) Generate $z_{ijk}^{*(b)} \sim \text{Bern}(p_{ijk}^{*(b)})$. If $z_{ijk}^{*(b)} = 1$, do $y_{ijk}^{*(b)} = 0$. If $z_{ijk}^{*(b)} = 0$, generate $y_{ijk}^{*(b)} \sim \text{NB}(r, \mu_{ijk}^{*(b)})$.
 - (c) Use $(y_{ijk}^{*(b)}, \mathbf{x}_{ijk})$ to calculate the ML estimate $\hat{\theta}_\ell^{*(b)}$.
- 3 Sort the values $\hat{R}_\ell^{*(b)} = D^{1/2} (\hat{\theta}_\ell^{*(b)} - \hat{\theta}_\ell)$, $b = 1, \dots, B$, from smallest to largest: $\hat{R}_{\ell(1)}^* \leq \dots \leq \hat{R}_{\ell(B)}^*$. A $(1 - \alpha)\%$ basic percentile bootstrap confidence interval for θ_ℓ is

$$(\hat{\theta}_\ell - D^{-1/2} \hat{R}_{\ell(\lfloor (\alpha/2)B \rfloor)}^*, \hat{\theta}_\ell + D^{-1/2} \hat{R}_{\ell(\lfloor (1-\alpha/2)B \rfloor)}^*).$$

MSE estimators and prediction intervals

- 1 Calculate the ML estimate $\hat{\theta}$. Run **Algorithm B2**.
- 2 Repeat B times ($b = 1, \dots, B$):
 - 1 Run steps (a) and (b) of Algorithm B1.
 - 2 For $i \in \mathbb{I}, j \in \mathbb{J}, k \in \mathbb{K}$, calculate $\mu_{y_{ijk}}^{*(b)} = (1 - p_{ijk}^{*(b)})\mu_{ijk}^{*(b)}$.
 - 3 Use the bootstrap sample $(y_{ijk}^{*(b)}, \mathbf{x}_{ijk})$, $i \in \mathbb{I}, j \in \mathbb{J}, k \in \mathbb{K}$, to calculate the ML estimate $\hat{\theta}^{*(b)}$ and the predictor $\hat{\mu}_{y_{ijk}}^{*(b)}$.
- 3 For $i \in \mathbb{I}, j \in \mathbb{J}, k \in \mathbb{K}$, calculate

$$mse^*(\hat{\mu}_{y_{ijk}}) = \frac{1}{B} \sum_{b=1}^B (\hat{\mu}_{y_{ijk}}^{*(b)} - \mu_{y_{ijk}}^{*(b)})^2, rmse^*(\hat{\mu}_{y_{ijk}}) = (mse^*(\hat{\mu}_{y_{ijk}}))^{1/2},$$

$$rrmse^*(\hat{\mu}_{y_{ijk}}) = \frac{rmse^*(\hat{\mu}_{y_{ijk}})}{\hat{\mu}_{y_{ijk}}}, \bar{y}_{ijk}^* = \sum_{b=1}^B \frac{y_{ijk}^{*(b)}}{B}, \text{var}^*(y_{ijk}) = \frac{1}{B-1} \sum_{b=1}^B (y_{ijk}^{*(b)} - \bar{y}_{ijk}^*)^2,$$

$$PI_{ijk}^{\alpha} = \left(\hat{\mu}_{y_{ijk}} - z_{1-\alpha/2}(\text{var}^*(y_{ijk}))^{1/2}, \hat{\mu}_{y_{ijk}} + z_{1-\alpha/2}(\text{var}^*(y_{ijk}))^{1/2} \right).$$

Data and problem of interest

Model and predictors

Bootstrap-based inference

Model-based statistical analysis

Wildfire forecasting and error measurements

Conclusions

Appendix: ML-Laplace approximation algorithm

Bernoulli regression parameters

Selected model (2002-2014)

- The Bernoulli (BE) mixed model has $q_1 = 6$ covariables.
 - $x_{1,1} = \text{intercept}$, $x_{1,2} = \text{hr}$, $x_{1,3} = \text{np.300}$,
 - $x_{1,4} = \text{ta.max}$, $x_{1,5} = \text{year3.2}$, $x_{1,6} = \text{year3.3}$.

β	$\hat{\beta}$	SE	z-value	$P(> z)$	CI 95% (asyp.)	CI 95% (boot)
$\beta_{1,1}$	-14.353	2.716	-5.28	0.000	(-19.67, -9.03)	(-17.73, -11.21)
$\beta_{1,2}^{hr}$	0.192	0.028	6.93	0.000	(0.138, 0.246)	(0.158, 0.232)
$\beta_{1,3}^{np.300}$	0.143	0.048	3.00	0.003	(0.049, 0.236)	(0.070, 0.213)
$\beta_{1,4}^{ta.max}$	-0.132	0.040	-3.39	0.001	(-0.208, -0.056)	(-0.175, -0.094)
$\beta_{1,5}^{year3.2}$	-1.048	0.219	-4.79	0.000	(-1.477, -0.619)	(-1.367, -0.768)
$\beta_{1,6}^{year3.3}$	-0.754	0.271	-2.78	0.005	(-1.285, -0.223)	(-1.253, -0.311)

Table: ML regression parameters of the Bernoulli mixed model.

Bernoulli regression parameters

If the remaining variables are fixed, the **negative signs** of the ML estimates of the regression parameters show that an **increase of**

- maximum temperature recorded (*ta.max*),

contribute to **reduce** the number of absolute zeros.

On the other hand, an **increase of**

- humidity (*hr*),
- percentage of days with precipitations greater than 30mm (*np.300*).

increases the number of absolute zeros.

- The **negative signs** of *year3.2* and *year3.3* show that the number of absolute zeros **has decreased** from 2002-2006 to 2007-2012 and 2013-2016.

Bernoulli variance parameters

The estimates of the **standard deviation parameters** are

- $\hat{\phi}_{11} = 0.1035$ for the month random effects $u_{1,j}$,
- $\hat{\phi}_{12} = 1.8717$ for the province random effects $u_{1,k}$.

The **asymptotic and bootstrap 95% CIs** are

- ϕ_{11} : (0.0037, 2.9343) and $(1.0621 \cdot 10^{-9}, 0.2972)$.
- ϕ_{12} : (1.4243, 2.4595) and (1.3193, 2.3050).
- No interval contains the zero: $u_{1,j}$ and $u_{1,k}$ are needed for modelling the mixture between 0 and the NB counts.

NB regression parameters

Selected model (2002-2014)

- The NB mixed model has $q_2 = 18$ covariables.
 - $x_{2,1} = \text{intercept}$,
 - $x_{2,2} = e$, $x_{2,3} = hr$,
 - $x_{2,4} = n.llu$, $x_{2,5} = n.nie$,
 - $x_{2,6} = np.300$, $x_{2,7} = nt.00$,
 - $x_{2,8} = nw.55$, $x_{2,9} = nw.91$,
 - $x_{2,10} = q.mar$, $x_{2,11} = q.max$,
 - $x_{2,12} = q.min$, $x_{2,13} = ta.max$,
 - $x_{2,14} = ta.min$,
 - $x_{2,15} = tm.mes$,
 - $x_{2,16} = tm.min$,
 - $x_{2,17} = year3.2$,
 - $x_{2,18} = year3.3$.

β	$\hat{\beta}$	SE	z-value	$P(> z)$	CI 95% (asyp.)	CI 95% (boot)
$\beta_{2,1}$	-10.51	3.801	-2.764	0.0057	(-17.960, -3.058)	(-12.409, -8.227)
$\beta_{2,2}^e$	0.009	0.001	7.245	0.0000	(0.007, 0.012)	(0.009, 0.010)
$\beta_{2,3}^{hr}$	-0.045	0.003	-13.366	0.0000	(-0.051, -0.038)	(-0.046, -0.043)
$\beta_{2,4}^{nllu}$	-0.014	0.001	-14.018	0.0000	(-0.016, -0.012)	(-0.015, -0.014)
$\beta_{2,5}^{nnie}$	-0.032	0.003	-9.295	0.0000	(-0.039, -0.025)	(-0.034, -0.030)
$\beta_{2,6}^{np.300}$	-0.029	0.007	-4.070	0.0000	(-0.043, -0.015)	(-0.034, -0.024)
$\beta_{2,7}^{nt.00}$	0.016	0.001	14.686	0.0000	(0.014, 0.018)	(0.016, 0.017)
$\beta_{2,8}^{nw.55}$	0.018	0.002	10.821	0.0000	(0.015, 0.021)	(0.017, 0.019)
$\beta_{2,9}^{nw.91}$	-0.025	0.008	-2.986	0.0028	(-0.041, -0.008)	(-0.030, -0.020)
$\beta_{2,10}^{q.mar}$	0.146	0.038	3.873	0.0001	(0.072, 0.220)	(0.125, 0.165)
$\beta_{2,11}^{q.max}$	-0.057	0.019	-2.973	0.0029	(-0.095, -0.019)	(-0.068, -0.046)
$\beta_{2,12}^{q.min}$	0.049	0.019	2.601	0.0093	(0.012, 0.086)	(0.039, 0.060)
$\beta_{2,13}^{ta.max}$	0.053	0.006	8.933	0.0000	(0.041, 0.064)	(0.050, 0.056)
$\beta_{2,14}^{ta.min}$	0.025	0.007	3.407	0.0007	(0.011, 0.040)	(0.021, 0.030)
$\beta_{2,15}^{tm.mes}$	0.064	0.023	2.823	0.0048	(0.020, 0.109)	(0.055, 0.074)
$\beta_{2,16}^{tm.min}$	-0.150	0.021	-7.215	0.0000	(-0.191, -0.109)	(-0.158, -0.140)
$\beta_{2,17}^{year3.2}$	-0.180	0.021	-8.549	0.0000	(-0.222, -0.139)	(-0.189, -0.171)
$\beta_{2,18}^{year3.3}$	-0.272	0.030	-9.016	0.0000	(-0.331, -0.212)	(-0.285, -0.257)

Table: ML regression parameters of the NB mixed model.

NB regression parameters

If the remaining variables are fixed, the **negative signs** of the ML estimates of the regression parameters show that an **increase of**

- humidity (*hr*),
- number of general rainy days (*n.llu*),
- number of on snowy days (*n.nie*)
- number of days whose precipitation is higher than 300mm (*np.300*),
- number of days with wind speed greater than or equal to 91km/h (*nw.91*),
- maximum absolute pressure (*q.max*),
- average minimum temperature (*tm.min*)

contribute to **reduce** the number of wildfires.

- The **negative signs** of *year3.2* and *year3.3* show that the count of fires **has decreased** from 2002-2006 to 2007-2012 and 2013-2016.

NB regression parameters

On the other hand, an increase of

- the average vapor pressure (e),
- days with minimum temperatures below 0°C ($nt.00$),
- days with wind speed greater than or equal to 55km/h ($nw.55$),
- mean pressure at sea level ($q.mar$),
- maximum minimum pressure ($q.min$),
- absolute maximum temperature ($ta.max$),
- average temperature ($tm.mes$),
- absolute minimum temperatures ($ta.min$),

increases the number of forest fires.

NB variance and dispersion parameters

The estimates of the **standard deviation parameters** are

- $\hat{\phi}_{21} = 0.3426$ for the month random effects $u_{2,j}$,
- $\hat{\phi}_{22} = 1.1555$ for the province random effects $u_{2,k}$.

The **asymptotic and bootstrap 95% CIs** are

- ϕ_{21} : (0.2264, 0.5187) and (0.1765, 0.4511).
- ϕ_{22} : (0.9485, 1.4078) and (0.9140, 1.3703).
- No interval contains the zero: $u_{2,j}$ and $u_{2,k}$ are needed.

The estimate of the **dispersion parameter** $\gamma = 1/r$ is $\hat{\gamma} = 2.1514$, with an SE of 0.0214.

The **asymptotic and bootstrap 95% CIs** for γ are

- (2.0630, 2.2437) and (2.0069, 2.7547), respectively.

Residual analysis

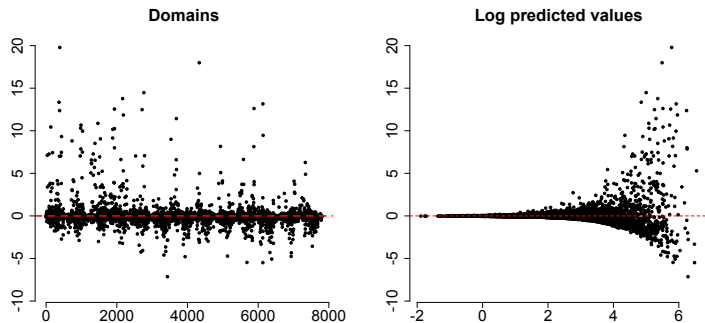


Figure: Standardized residuals vs domain indexes and log-predictions.

Residual analysis

Comments

- St.residuals fluctuate around zero, but there are more positive high residuals than negative ones.
- The asymmetry around zero is due to underprediction in provinces where the number of observed wildfires in summer is extremely high.
- The second plot also shows the mentioned asymmetry, which increases as the log-predictions grows.
- As the log-predictions increases, so does the variability of the residuals. This phenomenon is in agreement with the theoretical overdispersion of the BN distribution.

Residual analysis

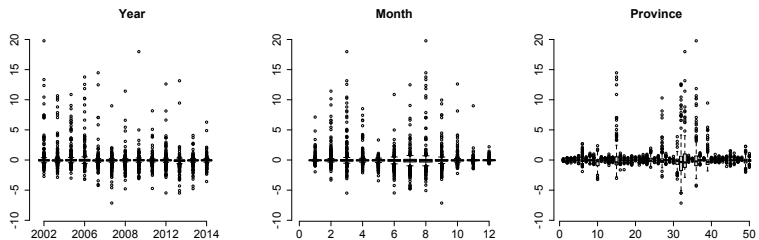


Figure: Boxplot of standardized residuals vs year, month and province.

Residual analysis

Comments

- Year does not seem to be discriminatory in the outliers's detection because no pattern is observed.
- Month seems to have more importance, with more atypical data at the end of winter and in summer.
- Province is a crucial factor in the outliers's detection.
- There are six provinces with absolute st.residuals greater than 3, with 82 observations (0.976% atypical data).
- The provinces are in the northwest of Spain: La Coruña (18), Lugo (5), Orense (22), Pontevedra (17), Asturias (14) and Cantabria (6).

Data and problem of interest

Model and predictors

Bootstrap-based inference

Model-based statistical analysis

Wildfire forecasting and error measurements

Conclusions

Appendix: ML-Laplace approximation algorithm

Forecasting

- This section deals with the prediction of the number of forest fires in 2015 (near future).
- As prediction scenario, we assume the 2015 data.
 - $\mathbf{x}_{1,ijk}$ and $\mathbf{x}_{2,ijk}$, $j \in \mathbb{J}$, $k \in \mathbb{K}$, $i = 2015$.
 - $(\beta'_1, \phi_{11}, \phi_{12})'$, $(\beta'_2, \phi_{21}, \phi_{22})'$, $\hat{\mathbf{u}}_{1,jk}$, $\hat{\mathbf{u}}_{2,jk}$, $j \in \mathbb{J}$, $k \in \mathbb{K}$.
- Based on the fitted AZINB11 model, we give predictions $\hat{\mu}_{y|ijk}$ for a short future horizon (one year).
- As the observed counts of wildfires y_{ijk} for 2015 are available, we can also check how good the predictions are.

Forecasting

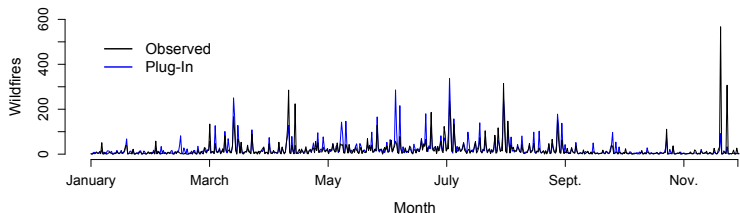
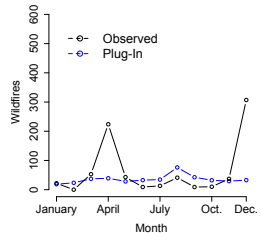
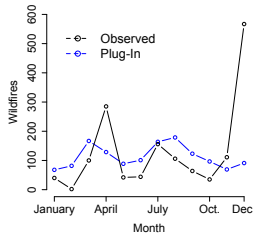
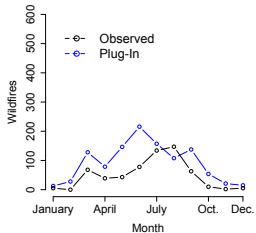
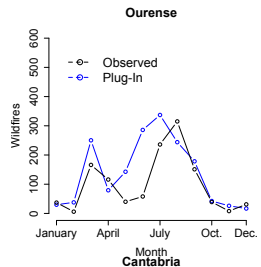
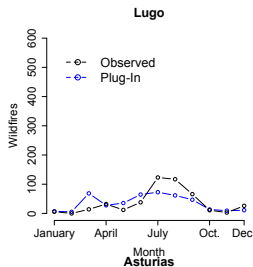
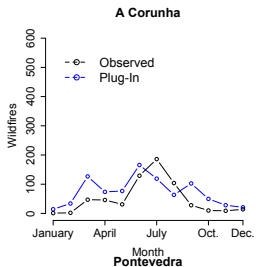


Figure: Wildfire forecasting for 2015 sorted by months and provinces (50 provinces per month).

- In general, predictions follows the trend of observations.
- Some predictions are far from observations in conflictive provinces (Asturias, Cantabria) during the summer. See next Figure.

Forecasting



Forecasting

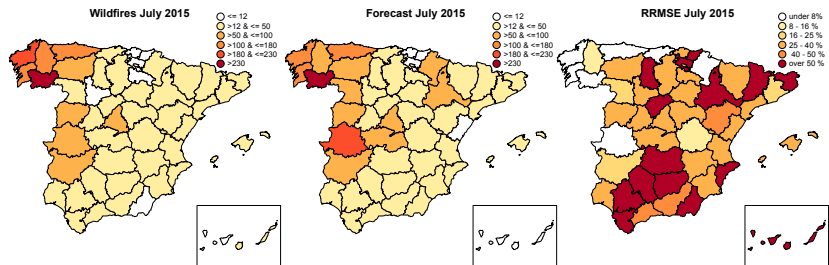


Figure: Observed (*left*) and predicted (*center*) wildfires and RRMSE (*right*) in July 2015.

Forecasting

Relative squared prediction errors (RSPE) for provinces in 2015

$$RSPE_{I,k} = \frac{\sqrt{\sum_{j=1}^{12} (y_{ljk} - \hat{\mu}_{ljk}^{in})^2}}{\sum_{j=1}^{12} y_{ljk}}, \quad I = 14$$

Coverage probabilities for provinces in 2015 ($I = 14$)

$$C_{I,k}^{\alpha} = \frac{1}{12} \sum_{j=1}^{12} C_{ljk}^{\alpha}, \quad C_{ljk}^{\alpha} = I(y_{ljk} \in PI_{ljk}^{\alpha}), \quad k = 1, \dots, K.$$

For months of 2015, the RSPEs and the coverage probabilities are

$$RSPE_{lj.} = \frac{\sqrt{\sum_{k=1}^{50} (y_{ljk} - \hat{\mu}_{ljk}^{in})^2}}{\sum_{k=1}^{50} y_{ljk}}, \quad C_{lj.}^{\alpha} = \frac{1}{50} \sum_{k=1}^K C_{ljk}^{\alpha}, \quad C_{ljk}^{\alpha} = I(y_{ljk} \in PI_{ljk}^{\alpha}),$$

Forecasting

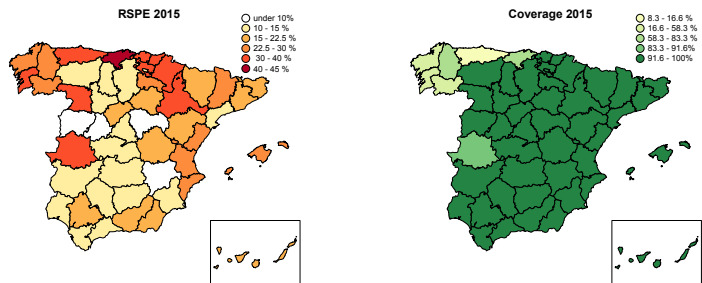


Figure: RSPE and coverage probabilities for each province in 2015, both in %.

Forecasting

	Jan.	Feb.	March	April	May	June
$RSPE_{lj}$.	16.96	20.30	14.41	20.57	18.43	24.32
C_{lj}^α .	94	88	88	90	92	94
	July	Aug.	Sept.	Oct.	Nov.	Dec.
$RSPE_{lj}$.	10.96	9.69	14.94	27.61	19.03	46.86
C_{lj}^α .	98	94	92	94	92	94

Table: Monthly RSPEs and coverage probabilities for 2015, both in %.

Data and problem of interest

Model and predictors

Bootstrap-based inference

Model-based statistical analysis

Wildfire forecasting and error measurements

Conclusions

Appendix: ML-Laplace approximation algorithm

Conclusions

- Zero-inflated negative binomial mixed models are flexible tools to
 - describe the behavior of the number of fires in Spanish provinces and months during the period 2002-2014.
 - to predict the number of fires in 2015.
- We calculated the ML estimators of the model parameters and the mode predictors of the random effects by applying the Laplace approximation algorithm.
- We introduced plug-in predictors of number of fires and parametric bootstrap estimators of their MSEs.
- For the out-of-sample data (2015), we also gave prediction intervals.

Data and R codes

- **Fire data:** General Forest Fire Statistics (EGIF) of the Spanish Government.
<https://www.miteco.gob.es/es/biodiversidad/servicios/banco-datos-naturaleza/informacion-disponible/incendios-forestales.html>
- **Meteorological data:** Spanish Meteorological Agency (AEMET)
https://www.aemet.es/es/datos_abiertos
- **Labour market data:** Spanish Statistical Office (INE).
<https://www.ine.es/prodyser/microdatos.htm>
- **R codes:** Downloadable from the repository
<https://github.com/mbugallo/aZINB11Fires>
- **R package glmmTMB:** To fit the aZINB11 model.

Thank you

Thank you
for your attention

M. Bugallo, M.D. Esteban, M.F. Marey-Pérez, D. Morales (2023).
Wild fire prediction using zero-inflated negative binomial mixed models: Application to Spain.
Journal of Environmental Management 328, 116788, 1-16.

Data and problem of interest

Model and predictors

Bootstrap-based inference

Model-based statistical analysis

Wildfire forecasting and error measurements

Conclusions

Appendix: ML-Laplace approximation algorithm

Appendix

- $h : \mathbb{R}^m \mapsto \mathbb{R}$ is a twice continuously differentiable function with a global maximum at the column vector \mathbf{x}_0 , i.e.
- $\dot{h}(\mathbf{x}_0) = \left. \frac{dh}{d\mathbf{x}} \right|_{\mathbf{x}=\mathbf{x}_0} = 0$ and $\ddot{h}(\mathbf{x}_0) = \left. \frac{d^2h}{d\mathbf{x}^2} \right|_{\mathbf{x}=\mathbf{x}_0}$ is negative definite.
- A Taylor series expansion of $h(\mathbf{x})$ around \mathbf{x}_0 yields to

$$\begin{aligned} h(\mathbf{x}) &\approx h(\mathbf{x}_0) + \dot{h}'(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)' \ddot{h}(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) \\ &= h(\mathbf{x}_0) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)' \ddot{h}(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0). \end{aligned}$$

- The multivariate Laplace approximation is

$$\begin{aligned} \int_{\mathbb{R}^m} e^{h(\mathbf{x})} d\mathbf{x} &\approx \int_{\mathbb{R}^m} e^{h(\mathbf{x}_0)} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \mathbf{x}_0)' (-\ddot{h}(\mathbf{x}_0))(\mathbf{x} - \mathbf{x}_0) \right\} d\mathbf{x} \\ &= (2\pi)^{m/2} |-\ddot{h}(\mathbf{x}_0)|^{-1/2} e^{h(\mathbf{x}_0)}, \end{aligned} \quad (1)$$

Appendix

The **likelihood** of the **AZINB11 model** is

$$P(\mathbf{y}; \boldsymbol{\theta}) = \int_{\mathbb{R}^{4JK}} P(\mathbf{y}|\mathbf{u}; \boldsymbol{\theta}) f_{\mathbf{u}}(\mathbf{u}) d\mathbf{u} = \int_{\mathbb{R}^{4JK}} \exp \{h(\mathbf{u}; \mathbf{y}, \boldsymbol{\theta})\} d\mathbf{u}, \quad (2)$$

where

$$\begin{aligned} h(\mathbf{u}; \mathbf{y}, \boldsymbol{\theta}) &= \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K \log P(y_{ijk} | \mathbf{u}_{jk}; \boldsymbol{\theta}) \\ &\quad - \frac{4JK}{2} \log 2\pi - \frac{1}{2} \sum_{j=1}^J \sum_{k=1}^K (u_{1,j}^2 + u_{1,k}^2 + u_{2,j}^2 + u_{2,k}^2). \end{aligned}$$

- To apply the **Laplace approximation** to the integral in (2), we maximize $h(\mathbf{u}; \mathbf{y}, \boldsymbol{\theta})$ in \mathbf{u} , given \mathbf{y} and $\boldsymbol{\theta}$.
- For simplicity, we write $h(\mathbf{u})$.

Appendix

- The **Newton-Raphson** updating equation is

$$\mathbf{u}^{(i+1)} = \mathbf{u}^{(i)} - \ddot{h}^{-1}(\mathbf{u}^{(i)}) \dot{h}(\mathbf{u}^{(i)}). \quad (3)$$

- Let \mathbf{u}° be the argument of maxima of the function $h(\mathbf{u})$.
- It holds $\dot{h}(\mathbf{u}^\circ) = 0$ and the matrix $\ddot{h}(\mathbf{u}^\circ)$ is negative definite.
- The **loglikelihood** of the **AZINB11 model** can be approximated by

$$\log P(\mathbf{y}; \boldsymbol{\theta},) \approx 2JK \log 2\pi + h(\mathbf{u}^\circ) - \frac{1}{2} \log |-\ddot{h}(\mathbf{u}^\circ)| \triangleq g(\boldsymbol{\theta}; \mathbf{y}, \mathbf{u}^\circ).$$

Appendix

- The following step is to maximize $g(\boldsymbol{\theta}; \mathbf{y}, \mathbf{u}^\circ)$ in $\boldsymbol{\theta} \in \Theta$.
- For simplicity, we write $g(\boldsymbol{\theta})$.
- Let us define $M = \dim(\Theta) = q_1 + q_2 + 4$.
- Let \dot{g} and \ddot{g} denote the $M \times 1$ vector and the $M \times M$ matrix of first and second order partial derivatives of $g(\boldsymbol{\theta})$.
- The **Newton-Raphson** updating equation is

$$\boldsymbol{\theta}^{(i+1)} = \boldsymbol{\theta}^{(i)} - \ddot{g}^{-1}(\boldsymbol{\theta}^{(i)}) \dot{g}(\boldsymbol{\theta}^{(i)}). \quad (4)$$

- The final **ML-Laplace approximation algorithm** combines (3) and (4).

Appendix

- 1 Set the initial values $i = 0$, $\varepsilon_1 > 0$, $\varepsilon_2 > 0$, $\varepsilon_3 > 0$, $\varepsilon_4 > 0$, $\boldsymbol{\theta}^{(0)}$, $\boldsymbol{\theta}^{(-1)} = \boldsymbol{\theta}^{(0)} + \mathbf{1}$, $\mathbf{u}^{(0)} = \mathbf{0}$, $\mathbf{u}^{(-1)} = \mathbf{1}$, where $\mathbf{0}$ and $\mathbf{1}$ are column vectors of zeros and ones respectively.
- 2 Until $\|\boldsymbol{\theta}^{(i)} - \boldsymbol{\theta}^{(i-1)}\|_2 < \varepsilon_1$, $\|\mathbf{u}^{(i)} - \mathbf{u}^{(i-1)}\|_2 < \varepsilon_2$, do
 - 1 Apply algorithm (3) with seeds $\mathbf{u}^{(i)}$, convergence tolerance ε_3 and $\boldsymbol{\theta} = \boldsymbol{\theta}^{(i)}$ fixed. Output: $\mathbf{u}^{(i+1)}$.
 - 2 Apply algorithm (4) with seed $\boldsymbol{\theta}^{(i)}$, convergence tolerance ε_4 and $\mathbf{u} = \mathbf{u}^{(i+1)}$ fixed. Output: $\boldsymbol{\theta}^{(i+1)}$.
 - 3 $i \leftarrow i + 1$.
- 3 Output: $\hat{\boldsymbol{\theta}} = \boldsymbol{\theta}^{(i)}$ and $\hat{\mathbf{u}} = \mathbf{u}^{(i)}$.