

Analysis of 2012–2014 Brazil Tuberculosis Data

Joshua Robert

Sen Souradeep

Zhang Zhaoyuan

Holste Lion

March 2023

1 Introduction

Problem Statement

Analysis of tuberculosis (TB) data originating from Brazil using Generalized Additive Models (GAMs). Brazil is divided into 557 administrative microregions, and the available data comprises counts of TB cases in each microregion for each of the years from 2012 to 2014.

2 Exploratory Analysis of Data and Problem

The TB data from Brazil includes 1,671 entries or samples with 14 columns of numeric data types that specify the characteristics of each sample. The columns are self-explanatory because they are called Indigenous, Illiteracy, Urbanisation, Density, Poverty, Poor Sanitation, Unemployment, Timeliness, Year, TB, Population, Region, lon, and lat. TB stands for tuberculosis, whereas lon and lat stand for longitude and latitude. The dataset has no missing values in the technical sense, but it contains some abnormalities, which increases the amount of pre-processing needed. The region is stored as a continuous variable despite being a factor variable. Nonetheless, changing it depends on the task at hand. Moreover, the collection includes coordinates that describe the precise geospatial locations of the micro-regions listed in the column. The next part gives a detailed exploration of the data.

Data Exploration

An in-depth analysis of the datasets reveals that the mean and median values for the indigenous population are low, but the maximum value is 50, which suggests that there are individual areas where the indigenous population is concentrated and that these areas may be areas of potential poverty and poor sanitation and should be areas where there are more cases of tuberculosis. The mean and median illiteracy rates are only 14 and 11, respectively, indicating that illiteracy is not widespread. However, the maximum value of 41 suggests that there are specific backward areas with significant populations lacking access to education, which also suggests that the area seems poor and has poor sanitation. There are still some places that are less urbanised, where there may be more occurrences of tuberculosis, but the mean, median, and minimum values for urbanisation are 70 and 22, respectively, suggesting that most areas are highly urbanised. Based on the population density data, most locations can fit one person in a room, but the highest value of 1.6 highlights the existence of some places with very high population densities, which sharply raise the rate of TB transmission. The distribution of the poverty data suggests that each district has different poverty levels, with only a limited number of districts where poverty is not a significant issue. Although the general results on inadequate sanitation are low, the maximum score of 58 indicates that some districts have poor sanitation and substantial disease risk. Although average unemployment rates are low, a maximum value of 20 indicates that some isolated regions may experience severe economic hardship, protracted social unrest, and potentially significant morbidity rates. With a minimum value of 0, notification timeliness data has a fairly wide range.

Data from the actual world follow a normal distribution. However, some of the columns indicate otherwise. Timeliness, Unemployment, and Urbanization are approximations of a normal distribution with few extremes, whilst the remainder is multimodal normal distributions. The above suggests that employing semi-parametric or non-parametric models to demonstrate the relationship between the target and predictors would be advantageous. The target variable in this study is a risk, defined as ‘TB/Population’, whereas the remaining variables are possible predictors. As demonstrated in the table below, most features in the dataset exhibit some connection. As some features are correlated, basic regression cannot be used because it would yield false results; rather, models that account for the connection can be used. It is vital to note that some characteristics are anticipated to have positive correlations (tuberculosis versus population, density versus poverty) and vice versa. Specifically, population density, poverty, health conditions, unemployment, and notification timeliness are likely high due to the high population density, the low economic share per capita, the high poverty rate, and the high jobless rate.

3 Model

We want to model the count of cases TB_i by actually modelling ρ_i using

$$TB_i \sim \text{Pois}(\lambda_i = z_i \rho_i) \quad TB_i \text{ indep.} \\ \log(\lambda_i) = \log(z_i) + \log(\rho_i)$$

where TB_i is the count of TB cases, z_i is the total population. Model $\log(\rho_i)$ as

$$\log(\rho_i) = \sum_{j=1}^8 f_j(x_{i,j}) \\ f(x_i) = \sum_{k=1}^q \beta_k b_k(x_i)$$

where $x_{i,j}$ is the j th covariate (out of 8 socio-economic covariates) for the i th instance/datum in the dataset and $f(\cdot)$ is a smooth function of said covariate. Hence, the model boils down to

$$TB_i \sim \text{Pois}(\lambda_i = z_i \rho_i) \quad TB_i \text{ indep.} \\ \log(\lambda_i) = \log(z_i) + \sum_{j=1}^8 \sum_{k=1}^q \beta_{j,k} b_{j,k}(x_{i,j})$$

Looking at the distribution of the residuals of the model, we can see that the data is clearly far too overdispersed to be modelled by a Poisson, which has a fixed dispersion parameter. Even with 80 knots per smooth term the model doesn't seem to have enough flexibility which may be another indicator that a Poisson model is unsuitable for such overdispersed data. We propose the conventional alternative to the Poisson - the Negative Binomial model. Doing so, leads to a drop in the AIC. So the model distribution is changed to Negative Binomial with the same parameterisation except for the feature that the count of TB cases is now Negative Binomial distributed with mean as described above. See Table ?? appendix for a showcase of different model configurations and their associated AIC.

$$TB_i \sim \text{NB}(\lambda_i, \sigma_i^2) \quad TB_i \text{ indep.} \\ \lambda_i = z_i \rho_i; \quad \sigma_i^2 = \lambda_i + \frac{\lambda_i^2}{k} \\ \log(\lambda_i) = \log(z_i) + \sum_{j=1}^8 \sum_{k=1}^q \beta_{j,k} b_{j,k}(x_{i,j})$$

where k is a dispersion parameter, later estimated by the `gam` function in R.

When having a look at the relationship between the squared residuals and the fitted values one sees that the relation is not exactly quadratic, but rather close to 0, which would reflect the relation between model variance and the expected value in a Gaussian Distribution Model (additional evidence is provided by the Residuals vs. Fitted plot). However, fitting a Gaussian model leads to very skewed residuals, indicating that the data is apparently not Gaussian. So the model distribution is changed to Negative Binomial with the same parameterization except for the feature that the count of TB cases is now Negative Binomial distributed with mean λ_i as described above.

Given this base model, we investigate whether all given socio-economic variables are needed to explain the response or whether there exists a model with fewer parameters. The p-value for the smooth term of Illiteracy points towards it not being statistically significant. Poverty, although not statistically insignificant, has the second-largest p-value. These terms are sequentially dropped and the resulting model checked against the original model via a Likelihood Ratio Test (conducted using the `anova` function in R). We find that leaving out Illiteracy does not alter the model at a 5%-level of significance, whereas taking out both Poverty and Illiteracy does. So, in the following, we use a model with all of the socio-economic variables except Illiteracy. Note that this converts our linear predictor to

$$\log(\lambda_i) = \log(z_i) + \sum_{j=1}^7 \sum_{k=1}^q \beta_{j,k} b_{j,k}(x_{i,j})$$

This leaves us with a model with $\text{AIC} = 14,391.19$ and 43.9% of deviance explained. Running `gam.check()` lets us analyse the residual plots (see Figure) and examine the basis functions for the model. The QQ plot tells us that the model fails to predict well on the upper and lower ends of the response variable. Increasing the knots to 20 per covariate leads to marginal improvement with 44.9% deviance explained. More efficient extensions can be to add 1) spatial, 2) temporal and 3) spatio-temporal covariates.

First, we will try adding spatial terms. The spatial model adds a smoothed term which is function of the longitude and the latitude. A bivariate function is used because it makes sense to assume that there are more cases at certain locations

(defined by the interaction between latitude and longitude) than others, rather than there being more cases at locations with a certain longitude for any latitude, or the other way round. Hence, our linear predictor is now

$$\log(\lambda_i) = \log(z_i) + \sum_{j=1}^7 \sum_{k=1}^q \beta_{j,k} b_{j,k}(x_{i,j}) + \sum_{k=1}^q \beta_k b_k(\text{lon}_i, \text{lat}_i)$$

Using this model with the regular **s** smoother function from the **mgcv** package leads to a model that can explain 56.4% of the deviance and has a slightly lower AIC of 14,013.13. The QQ plot still points to the upper and lower tails being incorrectly predicted. At the cost of significantly more computation, using a tensor product smooth **te** on the bivariate spatial term with 20 knots allows us to make a decent improvement on this. See Appendix for different numbers of knots that were tested. This gets us to 69.9% deviance explained. The QQ plot looks considerably better with only a few problematic instances at the top and bottom quantiles.

We contest this with an extension on the model with only socio-economic covariates, but instead of adding spatial terms, we add the temporal dimension **Year**. The linear predictor becomes

$$\log(\lambda_i) = \log(z_i) + \sum_{j=1}^7 f_{2012,j}(x_{i,j}) \times x_{2012} + \sum_{j=1}^7 f_{2013,j}(x_{i,j}) \times x_{2013} + \sum_{j=1}^7 f_{2014,j}(x_{i,j}) \times x_{2014}$$

where the new terms $x_{2012}, x_{2013}, x_{2014}$ are indicator variables equating to 1 if **Year** is respectively 2012, 2013, 2014 and zero otherwise. Exercising some shorthand, it can be expressed as

$$\log(\lambda_i) = \log(z_i) + \sum_{t=2012}^{2014} \sum_{j=1}^7 f_{t,j}(x_{i,j}) \times x_t$$

where x_t is now the indicator variable for **Year**. A slightly separate approach can be tested with **Year** as a covariate instead of a grouping variable. In that case, the linear predictor would be

$$\log(\lambda_i) = \log(z_i) + \sum_{j=1}^7 f_{t,j}(x_{i,j}) + \sum_{t=2012}^{2014} \beta_t x_t$$

Neither of the temporal formulations show much increase in deviance explained (the one with year as grouping variable actually shows a decrease to 41.5%!). Their QQ plots are also much worse than the spatial model, showing gross deviations on high as well as low quantiles. Finally, we create a spatio-temporal model, including both **Year** as well as **lon, lat**. Its linear predictor is formulated as below

$$\log(\lambda_i) = \log(z_i) + \sum_{t=2012}^{2014} \left(\sum_{j=1}^7 \sum_{k=1}^q \beta_{j,k} b_{j,k}(x_{i,j}) + \sum_{k=1}^q \beta_k b_k(\text{lon}_i, \text{lat}_i) \right) \times x_t$$

This is a model which includes the term for the location and estimates a functional relation for each year and each explaining variable. The AIC of this model does not drop compared to the spatial model, so the spatial model (given that it is simpler) is the model we choose to best explain the ratio of TB cases per capita. To recall, it is formulated as

$$\log(\lambda_i) = \log(z_i) + \sum_{j=1}^7 \sum_{k=1}^q \beta_{j,k} b_{j,k}(x_{i,j}) + \sum_{k=1}^q \beta_k b_k(\text{lon}_i, \text{lat}_i)$$

Let us now have a closer look at the fit of the spatial model: It fits well even though the largest residuals are higher than expected from the model distribution. For districts that have a high number of cases, the predictor does not seem as accurate. But the highest residuals do not arise when the ratio of TB cases per capita is extraordinarily high, but rather when the absolute number of TB cases is high (see residuals vs. response). The variance of the model still seems too low for those values given that there are some predicted values in that high segment of response values (absolute number of TB cases) where the prediction for the response value is lower than the actual value, and some where the prediction of the actual value is higher than the actual value. Using this model, we predict the rate of TB per 100,000 inhabitants. See Figure ??

4 Critical Review and Conclusion

Drawbacks of the Model:

1. Predictions do not cover full range of data, as evinced by deviations in the QQ plot
2. ??

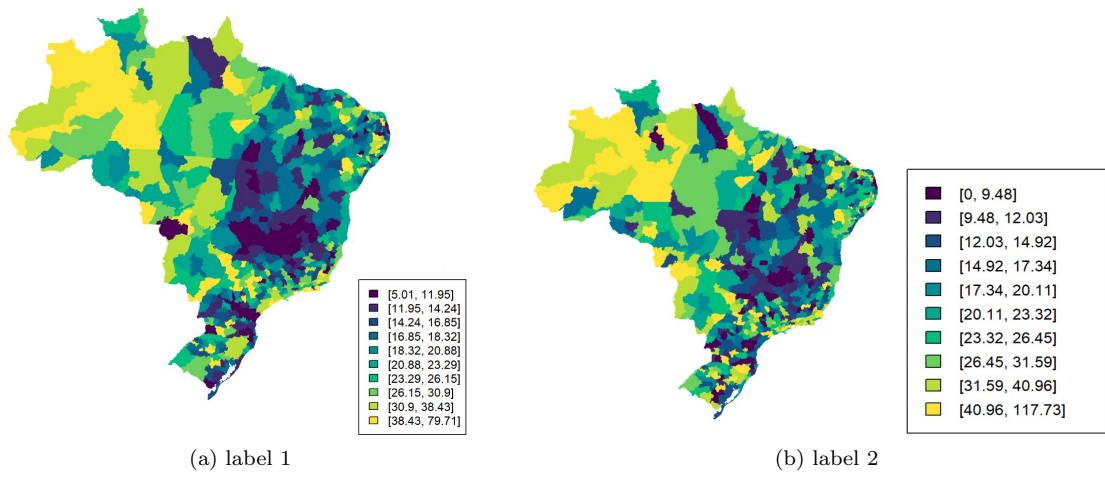


Figure 1: Predicted (a) and True(b) rates of TB per 100k inhabitants

5 References

Wood, S. N. (2017). Generalized Additive Models: An Introduction with R (2nd ed.). CRC Press.

6 Appendix

Figures

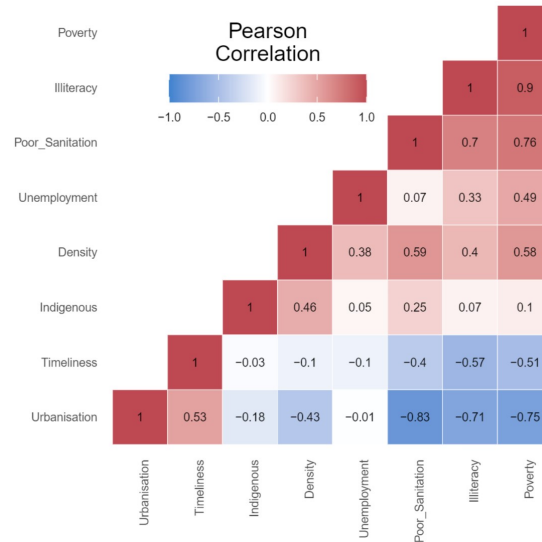


Figure 2: Correlogram shows covariates with highest positive and negative correlations.

Tables

Item	Quantity
Widgets	42
Gadgets	13

Table 1: An example table.

Code

```

1 #####SS CODE#####
2 # Import Libraries
3 library(mgcv) # required for GAM
4 library(tidyverse)
5 library(ggplot2) # required for plotting
6 library(dplyr) # required for filtering dataset
7 library(fields) # required for maps

```

```

8 library(maps) # required for maps
9 library(reshape2) # only required for melt in corr plot
10 library(car) # only required for VIF
11
12 # Load Data
13 load("C:/Users/soura/Documents/COMM511/group_coursework/datasets_project.RData")
14
15 # Investigate correlation
16 ### Resource -> http://www.sthda.com/english/wiki/ggplot2-quick-
17 ### correlation-matrix-heatmap-r-software-and-data-visualization
18 cormat <- cor(TBdata[,c(1,2,3,4,5,6,7,8)])
19 # Reorder
20 reorder_cormat <- function(cormat){
21   # Use correlation between variables as distance
22   dd <- as.dist((1-cormat)/2)
23   hc <- hclust(dd)
24   cormat <- cormat[hc$order, hc$order]
25 }
26
27 # Reorder the correlation matrix
28 cormat <- reorder_cormat(cormat)
29 # Get lower triangular matrix
30 cormat[lower.tri(cormat)] <- NA
31
32 melted_cormat <- melt(cormat , na.rm = TRUE)
33 melted_cormat$value = round(melted_cormat$value, 2)
34
35 # Create a ggheatmap
36 ggheatmap <- ggplot(melted_cormat, aes(Var2, Var1, fill = value))+
37   geom_tile(color = "white")+
38   scale_fill_gradient2(low = "#1a85d6", high = "#cf3e4f", mid = "white",
39     midpoint = 0, limit = c(-1,1), space = "Lab",
40     name="Pearson\nCorrelation") +
41   theme_minimal()+ # minimal theme
42   theme(axis.text.x = element_text(angle = 90, vjust = 1,
43     size = 12, hjust = 1))+
44   coord_fixed()
45
46 # Add correlation coefficients
47 ggheatmap +
48   geom_text(aes(Var2, Var1, label = value), color = "black", size = 2) +
49   theme(
50     axis.text.x = element_text(size = 6),
51     axis.text.y = element_text(size = 6),
52     axis.title.x = element_blank(),
53     axis.title.y = element_blank(),
54     panel.grid.major = element_blank(),
55     panel.border = element_blank(),
56     panel.background = element_blank(),
57     axis.ticks = element_blank(),
58     legend.justification = c(1, 0),
59     legend.position = c(0.6, 0.7),
60     legend.direction = "horizontal",
61     legend.text = element_text(size = 6)
62   ) +
63   guides(fill = guide_colorbar(barwidth = 7, barheight = 1,
64     title.position = "top", title.hjust = 0.5))
65
66 ##### Illiteracy is highly correlated with Poverty
67 ##### Carry out a Variance Inflation Test
68 model_all <- lm(TB ~ . , data = select(TBdata, 'Indigenous' , 'Illiteracy' ,
69   'Urbanisation' , 'Density' , 'Poverty' , 'Unemployment' , 'Timeliness' , 'Year' ,
70   'TB' , 'Population')) # with all the independent variables
71
72 vif(model_all) # Several variables are highly correlated
73
74 model_no_illiteracy <- lm(TB ~ . , data = select(TBdata, 'Indigenous',

```

```

75 'Urbanisation' , 'Density' , 'Poverty' , 'Unemployment' , 'Timeliness' , 'Year' ,
76 'TB' , 'Population')) # with all the independent variables
77
78 vif(model_no_illiteracy) # Poverty and Unemployment still seem highly correlated
79
80 model_no_illiteracy_no_poverty <- lm(TB ~ . , data = select(TBdata, 'Indigenous',
81 'Urbanisation' , 'Density', 'Unemployment' , 'Timeliness' , 'Year' ,
82 'TB' , 'Population')) # with all the independent variables
83
84 vif(model_no_illiteracy_no_poverty) # almost no variable is highly correlated
85
86 ## More formal tests are conducted to confirm the dropping of Illiteracy.
87 ## Check to see if Poverty should be dropped as well
88 prelim.model.1 <- gam(formula = TB ~ offset(log(Population)) + s(Indigenous) +
89 s(Illiteracy) + s(Urbanisation) + s(Density) + s(Poverty) + s(Poor_Sanitation)
90 + s(Unemployment) + s(Timeliness),
91 data = TBdata ,
92 family = nb(link = 'log')
93 )
94 # Show summary
95 summary(prelim.model.1)
96 par(mfrow=c(2,2))
97 gam.check(prelim.model.1)
98
99 ### Only the effect of illiteracy cannot be reliably stated to be non-zero
100 prelim.model.2 <- gam(formula = TB ~ offset(log(Population)) + s(Indigenous)
101 + s(Urbanisation) + s(Density) + s(Poverty) + s(Poor_Sanitation)
102 + s(Unemployment) + s(Timeliness),
103 data = TBdata ,
104 family = nb(link = 'log')
105 )
106 # Show summary
107 summary(prelim.model.2)
108
109 # Show summary
110 summary(prelim.model.2)
111 par(mfrow=c(2,2))
112 gam.check(prelim.model.2)
113
114 # Likelihood ratio test
115 anova(prelim.model.1 , prelim.model.2 , test = 'F') # p-value is over 0.05
116 # The models are statistically indistinguishable
117
118 ### Only the effect of illiteracy cannot be reliably stated to be non-zero
119 prelim.model.3 <- gam(formula = TB ~ offset(log(Population)) + s(Indigenous)
120 + s(Urbanisation) + s(Density) + s(Poor_Sanitation) + s(Unemployment)
121 + s(Timeliness),
122 data = TBdata ,
123 family = nb(link = 'log')
124 )
125 # Show summary
126 summary(prelim.model.3)
127
128 # Likelihood ratio test
129 anova(prelim.model.2 , prelim.model.3 , test = 'F') # p-value is less than 0.05
130 # The models are statistically different. Poverty should not be excluded.
131
132 ### Model chosen (with social covariates) is the negative binomial without
133 ### Illiteracy
134 summary(prelim.model.2) # Only 44% of the deviance is explained. Adding temporal
135 # and spatial covariates may improve this
136 par(mfrow=c(2,2))
137 gam.check(prelim.model.2)
138 par(mfrow = c(1,1))
139
140 ### Only the effect of illiteracy cannot be reliably stated to be non-zero
141 prelim.model.4 <- gam(formula = TB ~ offset(log(Population)) + s(Indigenous , k = 20)

```

```

142 + s(Urbanisation , k = 20) + s(Density , k = 20) + s(Poverty , k = 20)
143 + s(Poor_Sanitation , k = 20) + s(Unemployment , k = 20) + s(Timeliness , k = 20),
144 data = TBdata ,
145 family = nb(link = 'log')
146 )
147 # Show summary
148 summary(prelim.model.4)
149 par(mfrow=c(2,2))
150 gam.check(prelim.model.4)
151
152
153 ### Adding spatial covariates
154 spatial.model <- gam(formula = TB ~ offset(log(Population)) + s(Indigenous)
155 + s(Urbanisation) + s(Density) + s(Poor_Sanitation) + s(Unemployment) + s(Poverty)
156 + s(Timeliness) + s(lon , lat),
157 data = TBdata ,
158 family = nb(link = 'log')
159 )
160 # Check summary
161 summary(spatial.model)
162 # Check the smooth functions of the covars
163 plot(spatial.model)
164 par(mfrow=c(2,2))
165 gam.check(spatial.model)
166 par(mfrow = c(1,1))
167 spatial.model$aic
168
169 ### Using separate smoothers
170 spatial.model.2 <- gam(formula = TB ~ offset(log(Population)) + s(Indigenous)
171 + s(Urbanisation) + s(Density) + s(Poor_Sanitation) + s(Unemployment) + s(Poverty)
172 + s(Timeliness) + te(lon , lat , k = 20),
173 data = TBdata ,
174 family = nb(link = 'log')
175 )
176 # Check summary
177 summary(spatial.model.2)
178 # Check the smooth functions of the covariates
179 plot(spatial.model.2)
180 par(mfrow=c(2,2))
181 gam.check(spatial.model.2 , pch = 20)
182 par(mfrow = c(1,1))
183
184 ### PREDICTIONS
185 fitted_nb <- predict(spatial.model.2, newdata = TBdata , type = 'response')
186
187 # PLOTTING STUFF - ERROR WHILE COMPILING - WILL FIX LATER
188 par(mfrow = c(1,1))
189 plot.map(exp(log(fitted_nb) - log(TBdata$Population))*100000 , n.levels = 10)
190
191
192 TBdata$Year.asFactor <- factor(TBdata$Year)
193
194 #### Temporal covariates
195 temporal.model <- gam(formula = TB ~ offset(log(Population)) + s(Indigenous)
196 + s(Urbanisation) + s(Density) + s(Poor_Sanitation) + s(Unemployment) + s(Poverty)
197 + s(Timeliness) + Year.asFactor,
198 data = TBdata ,
199 family = nb(link = 'log')
200 )
201 # Check summary
202 summary(temporal.model) # Temporal alone doesn't add much to explaining the variance
203 # Check the smooth functions of the covariates
204 plot(temporal.model)
205 par(mfrow=c(2,2))
206 gam.check(temporal.model)
207 par(mfrow = c(1,1))
208

```

```

209 ##### Temporal covariates
210 temporal.model.2 <- gam(formula = TB ~ offset(log(Population))
211 + s(Indigenous , by=Year.asFactor) + s(Urbanisation , by=Year.asFactor)
212 + s(Density , by=Year.asFactor) + s(Poor_Sanitation , by=Year.asFactor)
213 + s(Unemployment , by=Year.asFactor) + s(Poverty , by=Year.asFactor)
214 + s(Timeliness , by=Year.asFactor) ,
215 data = TBdata ,
216 family = nb(link = 'log')
217 )
218 # Check summary
219 summary(temporal.model.2) # Temporal alone doesn't add much to explaining the variance
220 # Check the smooth functions of the covariates
221 plot(temporal.model.2)
222 par(mfrow=c(2,2))
223 gam.check(temporal.model.2)
224 par(mfrow = c(1,1))
225
226 ### Spatio-temporal model
227 spatio.temporal.model <- gam(formula = TB ~ offset(log(Population)) + s(Indigenous)
228 + s(Urbanisation) + s(Density) + s(Poor_Sanitation) + s(Unemployment) + s(Poverty)
229 + s(Timeliness) + te(lon , lat , k = 20) + Year.asFactor,
230 data = TBdata ,
231 family = nb(link = 'log')
232 )
233 # Check summary
234 summary(spatio.temporal.model) # Temporal alone doesn't add much to explaining the variance
235 # Check the smooth functions of the covariates
236 plot(spatio.temporal.model)
237 par(mfrow=c(2,2))
238 gam.check(spatio.temporal.model , pch = 20)
239 par(mfrow = c(1,1))
240
241 # PLOT
242 fitted_nb <- predict(spatio.temporal.model, newdata = TBdata , type = 'response')
243 TBdata$pred_rate <- fitted_nb/TBdata$Population*100000
244
245 # PLOTTING STUFF - ERROR WHILE COMPILING - WILL FIX LATER
246 par(mfrow = c(1,1))
247 plot.map(TBdata$pred_rate , n.levels = 10)
248
249 ### WIHTOUT POVERTY
250 ##### Spatio-temporal model - Poverty
251 spatio.temporal.model.wo.pov <- gam(formula = TB ~ offset(log(Population)) + s(Indigenous)
252 + s(Urbanisation) + s(Density) + s(Poor_Sanitation) + s(Unemployment)
253 + s(Timeliness) + te(lon , lat , k = 20) + Year.asFactor,
254 data = TBdata ,
255 family = nb(link = 'log')
256 )
257
258 # Check summary
259 summary(spatio.temporal.model.wo.pov) # Temporal alone doesn't add much to explaining the variance
260 # Check the smooth functions of the covariates
261 plot(spatio.temporal.model.wo.pov)
262 par(mfrow=c(2,2))
263 gam.check(spatio.temporal.model.wo.pov , pch = 20)
264 par(mfrow = c(1,1))
265
266 # Anova test
267 anova(spatio.temporal.model.wo.pov , spatio.temporal.model , test = 'F')
268
269 ### Without Indigenous
270 spatio.temporal.model.wo.indig <- gam(formula = TB ~ offset(log(Population))
271 + s(Urbanisation) + s(Density) + s(Poor_Sanitation) + s(Unemployment) + s(Poverty)
272 + s(Timeliness) + te(lon , lat , k = 20) + Year.asFactor,
273 data = TBdata ,
274 family = nb(link = 'log')
275 )

```



```

276
277 # Check summary
278 summary(spatio.temporal.model.wo.indig) # Temporal alone doesn't add much to explaining the variance
279 # Check the smooth functions of the covariates
280 plot(spatio.temporal.model.wo.indig)
281 par(mfrow=c(2,2))
282 gam.check(spatio.temporal.model.wo.indig , pch = 20)
283 par(mfrow = c(1,1))
284
285 # Anova test
286 anova(spatio.temporal.model.wo.indig , spatio.temporal.model , test = 'F')
287 # NOT ENOUGH EVIDENCE TO REMOVE INDIGENOUS
288
289 ## Spatio-temporal model
290 spatio.temporal.model.true <- gam(formula = TB ~ offset(log(Population)) + s(Indigenous)
291 + s(Urbanisation) + s(Density) + s(Poor_Sanitation) + s(Unemployment) + s(Poverty)
292 + s(Timeliness) + te(lon , lat , Year , k = 3),
293 data = TBdata ,
294 family = nb(link = 'log')
295 )
296
297 # Check summary
298 summary(spatio.temporal.model.true) # Temporal alone doesn't add much to explaining the variance
299 # Check the smooth functions of the covariates
300 plot(spatio.temporal.model.true)
301 par(mfrow=c(2,2))
302 gam.check(spatio.temporal.model.true , pch = 20)
303 par(mfrow = c(1,1))
304
305 spatio.temporal.model$aic
306 spatio.temporal.model.true$aic
307 spatio.temporal.model.wo.indig$aic
308
309
310
311 ### Spatio-temporal model - no poor sanitation , indigenous
312 spatio.temporal.model.wo.indig.poor_s <- gam(formula = TB ~ offset(log(Population))
313 + s(Urbanisation) + s(Density) + s(Unemployment) + s(Poverty)
314 + s(Timeliness) + te(lon , lat , k = 20) + Year.asFactor,
315 data = TBdata ,
316 family = nb(link = 'log')
317 )
318 # Check summary
319 summary(spatio.temporal.model.wo.indig.poor_s) # Temporal alone doesn't add much to explaining the variance
320 # Check the smooth functions of the covariates
321 plot(spatio.temporal.model.wo.indig.poor_s)
322 par(mfrow=c(2,2))
323 gam.check(spatio.temporal.model.wo.indig.poor_s , pch = 20)
324 par(mfrow = c(1,1))
325
326 anova(spatio.temporal.model , spatio.temporal.model.wo.indig.poor_s , spatio.temporal.model.wo.indig , test = 'F')
327
328 #####LH CODE#####
329 library (mgcv)
330
331
332 par(mfrow = c(2,2))
333 #fit poisson model with socio-economic variables
334 model_poisson <- gam(formula = TB ~ offset(log(Population)) + s(Indigenous) + s(Illiteracy) + s(Urbanisation) + s(Density) + s(Poverty)
335 summary(model_poisson)
336 model_poisson$aic
337 par(mfrow = c(2,2),pch = 20)
338 gam.check(model_poisson)
339 # add flexibility
340 model_poisson <- gam(formula = TB ~ offset(log(Population)) + s(Indigenous, k = 80) + s(Illiteracy, k = 80) + s(Urbanisation, k = 80) + s(Density, k = 80) + s(Poverty, k = 80)
341 gam.check(model_poisson)
342 summary(model_poisson) # SIGNS OF OVERFIT

```

```

343 plot(model_poisson) # SIGNS OF OVERFIT
344 #fit negative binomial model with socioeconomic
345 model_nb <- gam(formula = TB ~ offset(log(Population)) + s(Indigenous) + s(Illiteracy) + s(Urbanisation) + s(Density) + s(Poverty)
346 summary(model_nb)
347 model_nb$aic
348 #fit a linear relation between squared residuals and prediction to see whether another model describes the variance-fitted values re
349 summary(lm(log(model_nb$residuals^2) ~ log(predict(model_nb, type = 'response'))))
350 #drop Illiteracy
351 model_nb_2 <- gam(formula = TB ~ offset(log(Population)) + s(Indigenous) + s(Urbanisation) + s(Density) + s(Poverty) + s(Poor_Sani
352
353
354 #LRT
355 anova.gam(model_nb_2, model_nb, test = 'LRT')
356 #Null hypothesis not rejected -> drop poverty
357 model_nb_3 <- gam(formula = TB ~ offset(log(Population)) + s(Indigenous) + s(Urbanisation) + s(Density) + s(Poor_Sanitation) + s(U
358 #LRT
359 anova.gam(model_nb_3, model_nb_2, test = 'LRT')
360 # NULL hypothesis rejected -> drop only illiteracy, not poverty
361
362 model_nb_final <- model_nb_2
363 summary(model_nb_final)
364 gam.check(model_nb_final)
365 fitted_nb <- predict(model_nb_final, type = 'response')
366
367 # PLOTTING STUFF - ERROR WHILE COMPILING - WILL FIX LATER
368 par(mfrow = c(1,2))
369 plot.map(log(fitted_nb) - log(TBdata$Population) , n.levels = 10)
370 plot.map(log(TBdata$TB) - log(TBdata$Population))
371
372 #temporal model
373 par(mfrow = c(2,2), pch = 20)
374 model_nb_time <- gam(formula = TB ~ offset(log(Population)) + s(Indigenous, by = Year) + s(Urbanisation, by = Year) + s(Density, by
375 summary(model_nb)
376 model_nb$aic
377
378 #spatial model
379 model_nb_space <- gam(formula = TB ~ offset(log(Population)) + s(Indigenous) + s(Urbanisation) + s(Density) + s(Poverty) + s(Poor_S
380 summary(model_nb_space)
381 model_nb_space$aic
382 gam.check(model_nb_space)
383 anova.gam(model_nb_space, model_nb_final, test = 'LRT')
384
385 #spatio-temporal model
386 model_nb_time_and_space <- gam(formula = TB ~ offset(log(Population)) + s(Urbanisation, by = Year) + s(Density, by = Year) + s(Pove
387 summary(model_nb_time_and_space)
388 model_nb_time_and_space$aic
389 gam.check(model_nb_time_and_space)
390
391 # PLOT
392 fitted_nb <- predict(model_nb_time_and_space, newdata = TBdata , type = 'response')
393
394 # PLOTTING STUFF - ERROR WHILE COMPILING - WILL FIX LATER
395 par(mfrow = c(1,1))
396 plot.map(exp(log(fitted_nb) - log(TBdata$Population))*100000 , n.levels = 10)

```