

# Exploratory analysis of accesses to support centers for gender-based violence in Apulia

## Data

The dataset employed regards the counts of accesses to gender-based violence support centers in the Apulia region by residence municipality of the women victims of violence during 2022. R codes to generate the dataset are in the R script posted here which this report is based on.

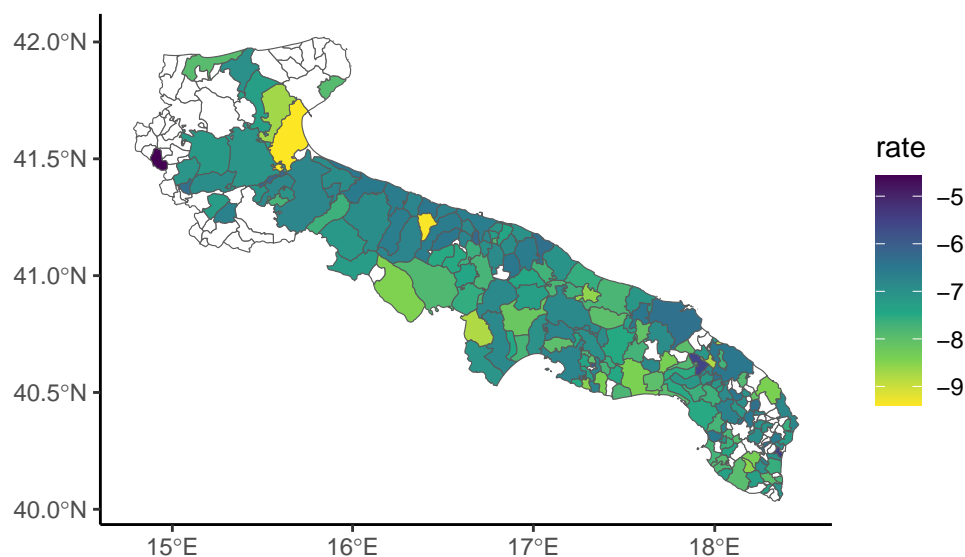
Here, we only take into account the violence reports which support centers actually take charge of, at the risk of underestimating the counts of gender-based violence cases. This choice is driven by the need of avoiding duplicated records, since e.g. it may happen that a support center redirects a victim to another support center.

In order to avoid singletons in the spatial structure of the dataset, we removed the Tremiti Islands from the list of municipalities included (0 accesses to support centers in 2022).

Therefore, the municipality-level dataset in scope consists of 256 observations.

We can only take into account the accesses to support centers for which the origin municipality of victims is reported. Therefore, the total count of accesses in scope is 2259. Among these accesses, 1516 were taken charge of.

Here, we plot the log-access rate per residence municipality, i.e. the logarithm of the ratio between access counts and female population. Blank areas correspond to municipalities from which zero women accessed support centers (82 municipalities).



## Covariates

Our target is explaining the number of accesses to support centers,  $y$ , defined at the municipality level, on the basis of a set of candidate known variables.

We model  $y$  via a simple Poisson GLM.

We have at disposal a number of candidate explanatory variables, which include the distance of a municipality from the closest support center and a set of variables measuring social vulnerability under different dimensions; these latter covariates are provided by the ISTAT. A more detailed description of these covariates is in this excel metadata file.

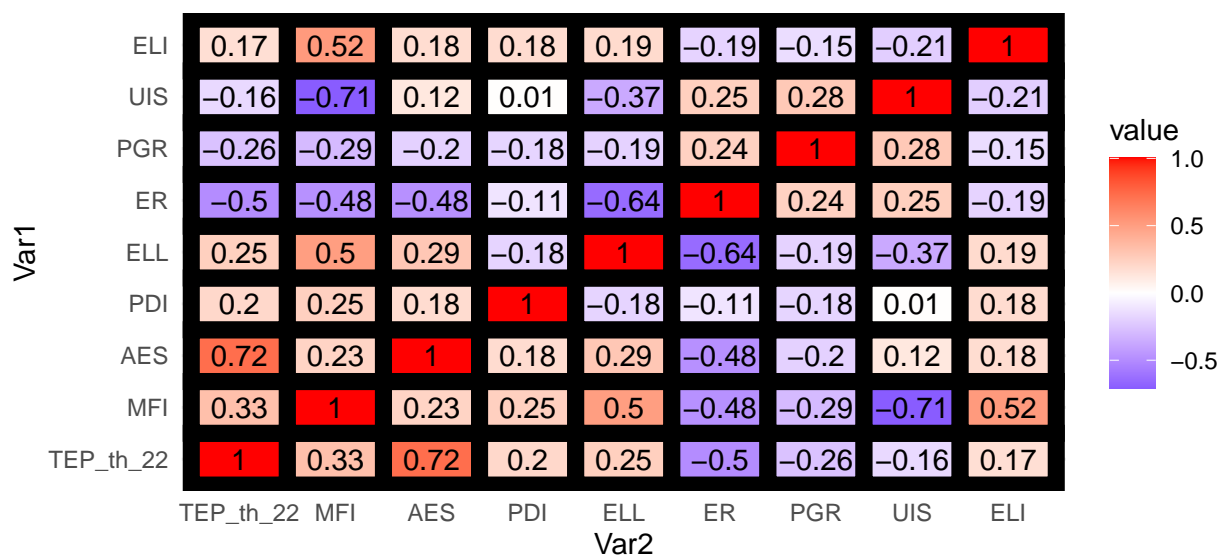
All covariates are scaled to have null mean and unit variance.

- TEP, i.e. the distance of each municipality from the closest municipality hosting a support center. Distance is measured by road travel time in minutes (acronym TEP stays for Tempo Effettivo di Percorrenza, i.e. Actual Travel Time).

For instance, the support center designated for the municipality of Adelfia (province of Bari, 3rd municipality in the dataset) is located in Capurso (BA). Then,  $TEP_3$  denotes the travel time between Adelfia and Capurso (17 minutes).

- AES, the distance from the closest infrastructural pole, always measured in travel time.
- MFI, i.e. the decile of municipality vulnerability index.
- PDI, i.e. the dependency index, i.e. population either  $\leq 20$  or  $\geq 65$  years over population in  $[20 - 64]$  years.
- ELL, i.e. the proportion of people aged  $[25 - 54]$  with low education.
- ERR, i.e. employment rate among people aged  $[20 - 64]$ .
- PGR, i.e. population growth rate with respect to 2011.
- UIS, i.e. the ventile of the density of local units of industry and services (where density is defined as the ratio between the counts of industrial units and population).
- ELI, i.e. the ventile of employees in low productivity local units by sector for industry and services.

First, we visualise the correlations among these explanatory variables:



We see the correlation between the two distances is very high (0.72), and so is the correlation between the fragility index decile and the density of productive units.

In the first case, we drop the distance from the nearest infrastructural pole. We do so because, if taken alone, the distance from the closest support center appears a slightly better predictor, using the Schwarz information criterion (or, indifferently, the Akaike Information Criterion):

```
stats::BIC(glm(N_ACC ~ 1 + AES, family = "poisson",
              offset = log(nn), data = dd_con))
```

```
## [1] 1124.922
```

```
stats::BIC(glm(N_ACC ~ 1 + TEP_th_22, family = "poisson",
               offset = log(nn), data = dd_con))
```

```
## [1] 1120.389
```

We should do the same for the other couple of variables but since MFI is a combination of all covariate except for TEP\_th, we will drop the synthetic indicator and leave the remainder.

## Nonspatial regression

We regress the counts of accesses  $y$  to support centers on the aforementioned explanatory variables. To estimate regression coefficients, all covariates are scaled to zero mean and unit variance.

$$y_i | \eta_i \sim \text{Poisson}(e^{\eta_i + P_i}) \quad \text{where} \quad \eta_i = X_i^\top \beta \quad (1)$$

Where  $X$  are the covariate defined earlier,  $\beta$  are covariate effects, and  $P_i$  is the female population aged  $\geq 15$  in municipality  $i$ .

To gain more insight on the role of all explanatory variables we show the posterior summaries of the full regression model

```
cav_glm <- glm(N_ACC ~ 1 + TEP_th_22 + ELI + PGR +
               UIS + ELL + PDI + ER,
               family = "poisson", offset = log(nn), data = dd_con)
summary(cav_glm)
```

```
##
## Call:
## glm(formula = N_ACC ~ 1 + TEP_th_22 + ELI + PGR + UIS + ELL +
##      PDI + ER, family = "poisson", data = dd_con, offset = log(nn))
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -7.46749    0.04525 -165.032 < 2e-16 ***
## TEP_th_22   -0.39714    0.04202  -9.451 < 2e-16 ***
## ELI         -0.07087    0.03369  -2.103  0.03543 *
## PGR          0.12915    0.04270   3.024  0.00249 **
## UIS         -0.09776    0.03511  -2.785  0.00536 **
## ELL         -0.11515    0.04427  -2.601  0.00930 **
## PDI         -0.06869    0.04522  -1.519  0.12872
## ER          -0.06641    0.05020  -1.323  0.18586
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
##
##      Null deviance: 729.12  on 255  degrees of freedom
## Residual deviance: 501.22  on 248  degrees of freedom
## AIC: 1083.8
##
## Number of Fisher Scoring iterations: 5
```

- TEP\_th\_22: The distance from the closest support center seems to play the key role. The easiest interpretation is that the physical distance represents a barrier to violence reporting. This is quite

intuitive if we think of the material dynamics of reporting gender-based violence: one could reasonably expect violent men to prevent their partners to come out and report the violence suffered.

- ELI: The (ventile of the distribution of the) share of employees in low productivity economic units is a clear indicator of (relative) economic underdevelopment. The most naive interpretation would be that in underdeveloped areas reporting gender violence is somewhat harder than in developed ones.
- PGR: The association with population growth rate is harder to interpret. This association is most likely influenced by several demographic instrumental variables we are not keeping into account and would indeed deserve a more dedicated focus.
- UIS: The (ventile of the distribution of the) density of production units has a somewhat ambiguous interpretation. From the one side, it has a strong negative relationship with the social frailty index. It should be therefore considered an indicator of economic development. Nevertheless, the regression coefficient bears the same sign as the incidence of low-productivity economic units. *Honestly I have no idea on how to interpret it.*
- ELL: The association with the proportion of people with low educational level has negative sign. The interpretation seems quite easy: cultural development, in general, would encourage reporting violence.
- PDI: The association with population dependency index does not seem significantly different from zero
- ER: Nor does the association with employment rate.

How do we interpret the regression coefficients? Keeping in mind we are working on the logarithm of the access rate, the standard deviation of the distance, expressed in minutes, is:

```
# Distance from closest support center  
attr(scale(dists_th_22$TEP_th_22), "scaled:scale")
```

```
## [1] 14.10021
```

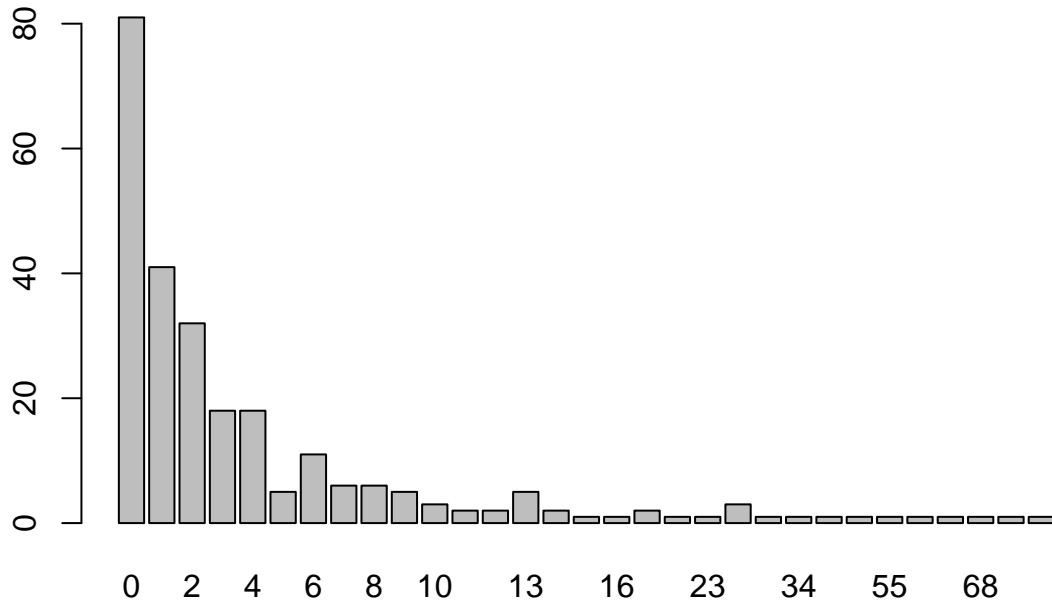
Hence e.g. each 14'6'' of distance of the a given municipality from the closest support center are associated with a decrease of 0.399 units in the log-frequency at which women from that municipality access to support centers.

Additionally, the number of zero counts is high:

```
sum(dd_con$N_ACC == 0)
```

```
## [1] 81
```

```
barplot(table(dd_con$N_ACC))
```



We may wonder if the data generating process incorporates a zero-generating component. We can model this augmented process through zero-inflated Poisson likelihood:

$$p(y_i|\eta_i) = \pi_0 \mathbb{I}\{y_i = 0\} + (1 - \pi_0) \frac{e^{-e^{\eta_i + P_i}} (e^{\eta_i + P_i})^{y_i}}{y_i!}$$

Where  $\pi_0 := \text{Prob}\{y_i = 0\}$  for all  $i$ . For the time being, we do not seek for explanatory variables for  $\pi_0$ . When explicitly accounting for the zero-generating process, we find the association of  $y$  with some explanatory variables to be reduced:

```
cav_zip <- pscl::zeroinfl(N_ACC ~ 1 + TEP_th_22 + ELI + PGR +
                        UIS + ELL + PDI + ER | 1, dist = "poisson",
                        link = "log", offset = log(nn), data = dd_con)

summary(cav_zip)$coefficients$count
```

	Estimate	Std. Error	z value	Pr(> z )
## (Intercept)	-7.40438222	0.04719042	-156.904359	0.000000e+00
## TEP_th_22	-0.40549864	0.04225844	-9.595683	8.336295e-22
## ELI	-0.04759633	0.03439110	-1.383972	1.663669e-01
## PGR	0.08421710	0.04464342	1.886439	5.923579e-02
## UIS	-0.06182367	0.03625139	-1.705415	8.811704e-02
## ELL	-0.11276640	0.04425341	-2.548197	1.082814e-02
## PDI	-0.05097672	0.04690631	-1.086778	2.771351e-01
## ER	-0.08627180	0.05080106	-1.698229	8.946463e-02

However, the MLE estimator for  $\pi_0$  is low:

```
summary(cav_zip)$coefficients$zero
```

```
##           Estimate Std. Error   z value   Pr(>|z|)
## (Intercept) -2.67672    0.362882 -7.376283 1.6277e-13
```

## Spatial regression

We plot the log-residuals  $\varepsilon$  of the GLM regression model, defined as  $\varepsilon := \ln y_i - \ln P_i - \ln \hat{y}_i$  being  $\hat{y}_i$  the fitted value.

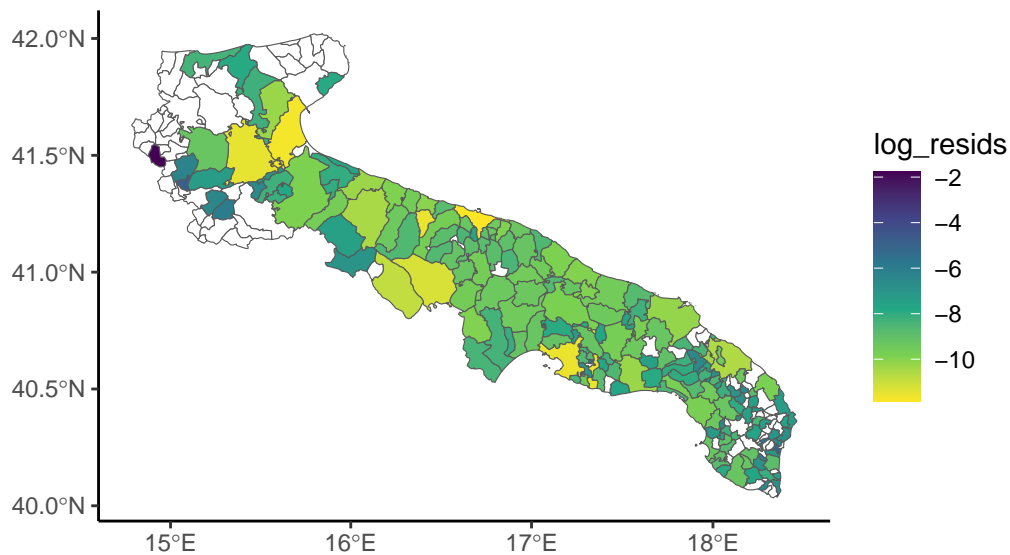


Figure 1: Log-residuals of glm regression using theoretical distance as explanatory variable

```
## Warning in spdep::poly2nb(dd_con[nonzero_con, ]): neighbour object has 2 sub-graphs;
## if this sub-graph count seems unexpected, try increasing the snap argument.
```

Residuals may exhibit spatial structure. To assess it, we employ the Moran and Geary tests. Since

Please notice that log-residuals only take finite values across the 175 municipalities whose female citizens have reported at least one case of violence in 2022.

Additionally, this set of municipalities includes 2 singletons, which we remove to assess the value of the Moran and Geary statistics. Thus, we have defined the indexes set `nonzero_con` as the set of municipalities from which at least one case of gender-based violence has been reported, *and* which have at least one neighbouring municipalities from which at least one case of gender-based violence was reported.

```
spdep::moran.test(resids_glm_th[nonzero_con],
                  listw = spdep::nb2listw(nb_con_nonzero))
```

```
##
## Moran I test under randomisation
##
## data:  resids_glm_th[nonzero_con]
## weights: spdep::nb2listw(nb_con_nonzero)
##
## Moran I statistic standard deviate = 5.9589, p-value = 1.269e-09
## alternative hypothesis: greater
## sample estimates:
```

```
## Moran I statistic      Expectation      Variance
##      0.316941893      -0.005813953      0.002933655

spdep::geary.test(resids_glm_th[nonzero_con],
                  listw = spdep::nb2listw(nb_con_nonzero))

##
## Geary C test under randomisation
##
## data:  resids_glm_th[nonzero_con]
## weights: spdep::nb2listw(nb_con_nonzero)
##
## Geary C statistic standard deviate = 4.3576, p-value = 6.576e-06
## alternative hypothesis: Expectation greater than statistic
## sample estimates:
## Geary C statistic      Expectation      Variance
##      0.731735532      1.000000000      0.003789984
```

In both cases, we find evidence for spatial autocorrelation. However, we must stress out this result does not refer to all the regional territory, but only to a subset of all municipalities (173 over 257)

Based on the autocorrelation evidence, though it has only been assessed for a subset of all municipalities, we try implementing some simple spatial models by adding a conditionally autoregressive latent effect, say  $z$ , to the linear predictor

$$\eta_i = X_i^\top \beta + z_i \quad (2)$$

We test a total of four models, all of which have a prior distribution depending on the spatial structure of the underlying graph, in this case the Apulia region.

We describe the spatial structure starting from municipality neighbourhood, and introduce the neighbourhood matrix  $W$ , whose generic element  $w_{ij}$  takes value 1 if municipalities  $i$  and  $j$  are neighbours and 0 otherwise. For each  $i \in [1, n]$ ,  $d_i := \sum_{j=1}^n w_{ij}$  is the number of neighbours of  $i$ -th municipality. Please notice we have  $n = 256$ .

For all models, we define  $\sigma^2$  as the scale parameter of the latent effect, and in order to avoid overfitting we set a PC-prior on it with rate parameter  $\lambda = 1.5$ , such that  $\text{Prob}(\sigma > \lambda) = 0.01$

Spatial models are computed by approximating the marginal posteriors of interest via the Integrated Nested Laplace Approximation (INLA), adopting the novel Variational Bayes Approach ?.

**ICAR model** The Intrinsic CAR model is the simplest formulation among spatial autoregressive models. The conditional distribution of each value  $z_i \mid z_{-i}$  is:

$$z_i \mid z_{-i} \sim N \left( \sum_{j=1}^n \frac{w_{ij}}{d_i} z_j, \frac{\sigma^2}{d_i} \right) \quad (3)$$

Since the joint distribution of  $z$  is improper, a sum-to-zero constraint is required for identifiability.

**PCAR model** The intrinsic autoregressive model is relatively simple to interpret and to implement, while also requiring the minimum number of additional parameter (either the scale or the precision).

The drawback, however, is that we implicitly assume a deterministic spatial autocorrelation coefficient equal to 1. When the autocorrelation is weak, setting an ICAR prior may be a form of misspecification.

A generalisation of this model is the PCAR (proper CAR), which introduces an autocorrelation parameter  $\alpha$ :

$$z_i \mid z_{-i} \sim N \left( \sum_{j=1}^n \alpha \frac{w_{ij}}{d_i} z_j, \frac{\sigma^2}{d_i} \right) \quad (4)$$

The R code to implement the PCAR model in R-INLA is in Appendix.

We show the posterior summary for the autocorrelation coefficient.

```
## Mean          0.812698
## Stdev         0.137317
## Quantile 0.025 0.456104
## Quantile 0.25  0.744349
## Quantile 0.5   0.849084
## Quantile 0.75  0.916071
## Quantile 0.975 0.974667
```

**BYM model** Perhaps, our data are generated by a process dominated from noise. We can thus try a different path: the BYM model. On a preliminary stance, we keep trusting in the accuracy of the Laplace approximation and stick to INLA. On a later stage, it would be more rigorous to compare INLA results to the posteriors of a model estimated with MCMC.

The BYM model we employ follows the parametrisation of (Riebler et al. 2016):

$$z_i = \sigma \left( \sqrt{\phi} u_i + \sqrt{1 - \phi} v_i \right) \quad (5)$$

where  $u$  is an ICAR field,  $v$  is an IID standard Gaussian white noise i.e.  $v \sim N(0, I)$ , and  $\phi$  is a mixing parameter  $\in [0, 1]$ .

**Leroux model** As an alternative to take into account both structured and unstructured latent effects, we also test the Leroux autoregressive model (Leroux, Lei, and Breslow 2000). In this case, the local prior for  $z_i$  is

$$z_i \mid z_{-i} \sim N \left( \sum_{j=1}^n \frac{\xi}{1 - \xi + \xi d_i} \frac{w_{ij}}{d_i} z_j, \frac{\sigma^2}{1 - \xi + \xi d_i} \right) \quad (6)$$

Where  $\xi \in [0, 1]$  is the mixing parameter. A more interesting representation of the Leroux model is the joint prior

$$z \mid \sigma^2, \xi \sim N(0, \sigma^2 [\xi R + (1 - \xi)I]^{-1})$$

where  $R := D - W$  is the graph Laplacian matrix,  $W$  is the neighbourhood matrix and  $D$  is the corresponding degree matrix.

**Comparison** We briefly compare these four through the WAIC (Gelman, Hwang, and Vehtari 2014):

```
## # A tibble: 5 x 5
##   Model  WAIC_Poisson WAIC_ZIP Eff_params_Poisson Eff_params_ZIP
##   <chr>      <dbl>    <dbl>          <dbl>          <dbl>
## 1 Null      1098.    1047.           20.4           18.1
## 2 ICAR      951.     963.           74.7           66.4
## 3 PCAR      947.     953.           76.8           65.0
## 4 BYM       943.     953.           75.3           65.1
## 5 Leroux    945.     952.           75.6           64.0
```



As we can see, models taking into account random noise have a better performance. Please notice the effective number of parameters, i.e. the number of *unconstrained* parameters is higher in the ICAR than in the BYM and Leroux models, even though these require an additional parameter.

## Model results

Here, posterior estimates of  $\beta$  under the BYM model are shown. We compare the results assuming both the Poisson and the ZIP likelihood:

##	Variable	Mean_Pois	Mean_ZIP	SD_Pois	SD_ZIP	q0.025_Pois	q0.025_ZIP
## 1	(Intercept)	-7.628	-7.533	0.066	0.069	-7.760	-7.671
## 2	TEP_th_22	-0.399	-0.433	0.078	0.072	-0.553	-0.575
## 3	ELI	-0.024	-0.028	0.064	0.058	-0.149	-0.142
## 4	PGR	0.138	0.108	0.073	0.068	-0.005	-0.026
## 5	UIS	-0.093	-0.081	0.068	0.061	-0.225	-0.199
## 6	ELL	-0.255	-0.215	0.085	0.078	-0.423	-0.371
## 7	PDI	0.006	-0.012	0.081	0.076	-0.154	-0.161
## 8	ER	-0.199	-0.180	0.101	0.092	-0.401	-0.365
##	q0.975_Pois	q0.975_ZIP					
## 1	-7.500	-7.400					
## 2	-0.245	-0.293					
## 3	0.102	0.087					
## 4	0.282	0.243					
## 5	0.041	0.039					
## 6	-0.090	-0.065					
## 7	0.165	0.138					
## 8	-0.003	-0.003					

Estimations of  $\beta$  differ from the nonspatial model. For all variables, credibility intervals are wider due to increased uncertainty.

- TEP\_th\_22 The effect of the distance from the closest support center remains similar in mean and the interpretation is not altered.
- ELI: The effect of the incidence of low-productivity economic units is shrunk in mean while its variability increases
- PGR: The association with population growth rate is still positive and significantly  $\neq 0$
- UIS: The association with the density of productive units appears not significant, due to increased variability
- ELL: The association with the incidence of low education levels, instead, is doubled in mean. We interpret this result as a strong *potential* impact of education on the chance that gender violence is reported
- PDI: The effect of structural dependency index is utterly negligible
- ER: The effect associated with employment rate is more than doubled in mean with respect to the nonspatial model. How to interpret this finding? Employment rate is clearly an indicator of economic development, hence the easiest interpretation is that - as it was with ELI under the nonspatial model - in more developed areas there is a higher chance that gender violence is reported.

Lastly, we take a look at model hyperparameters:

```
hyperpars <- tibble::tibble(
  Variable = c("Z_prec_Pois", "Mixing_Pois", "Zero_Prob", "Z_prec_ZIP", "Mixing_ZIP"),
  Mean = round(c(cav_bym_INLA$summary.hyperpar$mean, cav_bym_INLA_ZIP$summary.hyperpar$mean), 3),
  q0.025 = round(c(cav_bym_INLA$summary.hyperpar$q0.025quant, cav_bym_INLA_ZIP$summary.hyperpar$q0.025quant), 3),
  q0.975 = round(c(cav_bym_INLA$summary.hyperpar$q0.975quant, cav_bym_INLA_ZIP$summary.hyperpar$q0.975quant), 3)
```

```
Median = round(c(cav_bym_INLA$summary.hyperpar$`0.5quant`, cav_bym_INLA_ZIP$summary.hyperpar$`0.5quant`), 3)
q0.975 = round(c(cav_bym_INLA$summary.hyperpar$`0.975quant`, cav_bym_INLA_ZIP$summary.hyperpar$`0.975quant`), 3)
```

```
hyperpars
```

```
## # A tibble: 5 x 5
##   Variable      Mean q0.025 Median q0.975
##   <chr>         <dbl> <dbl> <dbl> <dbl>
## 1 Z_prec_Pois  2.80   1.36  2.63  5.22
## 2 Mixing_Pois  0.562  0.156 0.577  0.903
## 3 Zero_Prob    0.047  0.01  0.04  0.12
## 4 Z_prec_ZIP   5.74   1.95  5.02  13.9
## 5 Mixing_ZIP   0.373  0.013 0.317  0.924
```

Even though the value of the zero-inflation parameter is low, we can see deep differences in the structure of the latent effects: the precision parameter median of the ZIP model is almost double than under the Poisson model, while the mixing parameter is shrunk.

For completeness, we show the hyperparameters posterior summary also for the Leroux model. The mixing parameter is not directly comparable. Neither does the precision parameter, since only under intrinsic models can the Laplacian matrix be scaled *a priori*. That being said, we see the difference between the two likelihoods is analogous.

```
hyperpars_leroux <- tibble::tibble(
  Variable = c("Z_prec_Pois", "Mixing_Pois", "Zero_Prob", "Z_prec_ZIP", "Mixing_ZIP"),
  Mean = round(c(cav_leroux_INLA$summary.hyperpar$mean, cav_leroux_INLA_ZIP$summary.hyperpar$mean), 3),
  q0.025 = round(c(cav_leroux_INLA$summary.hyperpar$`0.025quant`, cav_leroux_INLA_ZIP$summary.hyperpar$`0.025quant`), 3),
  Median = round(c(cav_leroux_INLA$summary.hyperpar$`0.5quant`, cav_leroux_INLA_ZIP$summary.hyperpar$`0.5quant`), 3),
  q0.975 = round(c(cav_leroux_INLA$summary.hyperpar$`0.975quant`, cav_leroux_INLA_ZIP$summary.hyperpar$`0.975quant`), 3))
```

```
hyperpars_leroux
```

```
## # A tibble: 5 x 5
##   Variable      Mean q0.025 Median q0.975
##   <chr>         <dbl> <dbl> <dbl> <dbl>
## 1 Z_prec_Pois  1.62   0.856 1.54  2.81
## 2 Mixing_Pois  0.514  0.141 0.516  0.877
## 3 Zero_Prob    0.047  0.013 0.042  0.105
## 4 Z_prec_ZIP   2.86   1.21  2.64  5.81
## 5 Mixing_ZIP   0.409  0.079 0.39  0.826
```

It is, however, worth noticing how INLA manages to meet the regularity conditions to approximate the CPO more often than in the non-inflated model. Still, CPO must be re-computed manually before employing it as an evaluation metrics:

```
sum(cav_bym_INLA$cpo$failure)
```

```
## [1] 103.341
```

```
sum(cav_bym_INLA$cpo$failure > 0)
```

```
## [1] 223
```

```
sum(cav_bym_INLA_ZIP$cpo$failure)
```

```
## [1] 16.35669
```

```
sum(cav_bym_INLA_ZIP$cpo$failure > 0)
```

```
## [1] 111
```

```
sum(cav_leroux_INLA$cpo$failure)
```

```
## [1] 109.6474
```

```
sum(cav_leroux_INLA$cpo$failure > 0)
```

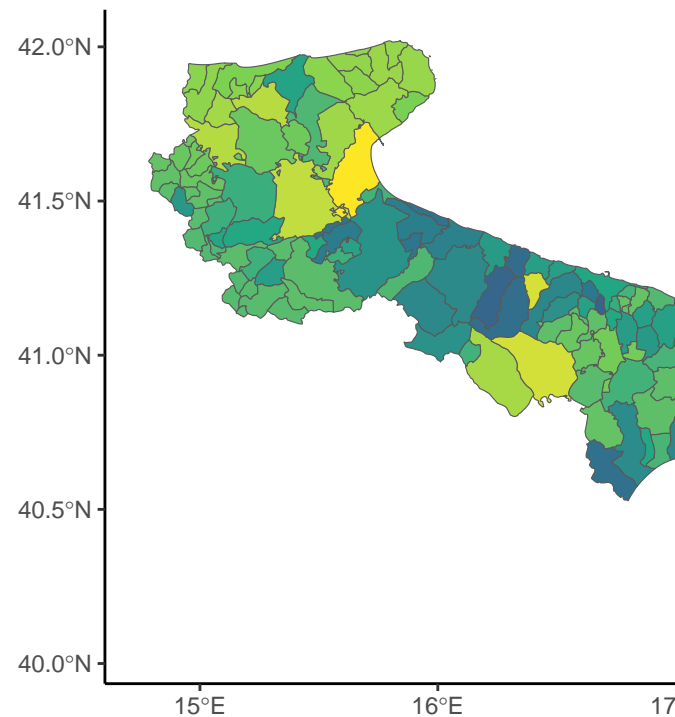
```
## [1] 221
```

```
sum(cav_leroux_INLA_ZIP$cpo$failure)
```

```
## [1] 22.50318
```

```
sum(cav_leroux_INLA_ZIP$cpo$failure > 0)
```

```
## [1] 87
```



Lastly, we plot the estimated latent BYM effect under the ZIP model

## Weakness elements and possible developments

From this preliminary analysis, inference on spatial models is hindered by the dominance of random noise over structured spatial effects. This can be argued from the posterior distribution of the mixing parameter in the BYM model, other than from the low spatial autocorrelation parameter in the PCAR.

This means that only to a small extent the variation in  $y$  not explained by covariates can be explained by spatial structure.

On the other hand, it is difficult to assert *all* variation not explained by covariates is pure noise, otherwise we

would have evidence for the lack of autocorrelation in residuals. We tested the hypothesis of no autocorrelation in GLM residuals by the Moran's  $I$  test, but in doing so we had to only test the residuals of areas with nonzero counts.

Moreover, spatial models are estimated using the INLA. While this is a broadly employed approach in epidemiology and in disease mapping, so far we did not assess how accurate the Laplace approximation has been.

To do so, we should e.g. rerun the same models using MCMC methods, e.g. using R libraries such as **CARBayes**, and replicating the same prior structure used.

Lastly, we did *not* model the rate at which gender violence occurs, but the occurrence of violence reports. Higher occurrence of violence reports from a given territory may thus depend on two factors: either the higher occurrence of violence in that territory, or the ease in reporting violence for the residents.

Whereas the easiest interpretation is that violence occurrence is underestimated in low-reporting areas, at the time being nothing prevents us from suspecting that the placement of support centers is at least partially strategic, i.e. the distribution of supporting centers is more dense in areas in which violence occurs, for some reason we don't know, as a higher frequency.

## Appendix: the WAIC

Following [?](#) , the WAIC is given by the sum of two components:

$$WAIC := 2 \sum_{i=1}^n \text{VAR}[\ln p(y_i|\theta)] - 2 \sum_{j=1}^n \ln \mathbb{E}[p(y_j|\theta)]$$

Where  $\theta$  is the full set of model parameters; the variance and the average are computed by integrating over the posterior of  $\theta$ . The first addendum denotes the number of free parameters, while the second term is a measure for goodness of fit.

## Appendix: R code to implement the PCAR model in INLA

Although it is not readily implemented in R-INLA (the "besagproper" effect is actually the Leroux model) we may base the R code on the 'INLAMSM' package (Palmí-Perales, Gómez-Rubio, and Martínez-Beneito 2021):

```
inla.rgeneric.PCAR.model <-
function (cmd = c("graph", "Q", "mu", "initial", "log.norm.const",
                  "log.prior", "quit"), theta = NULL) {
  interpret.theta <- function() {
    alpha <- 1/(1 + exp(-theta[1L])) # alpha modelled in logit scale
    mprec <- sapply(theta[2L], function(x) {
      exp(x)
    })
    PREC <- mprec
    return(list(alpha = alpha, mprec = mprec, PREC = PREC))
  }
  graph <- function() {
    G <- Matrix::Diagonal(nrow(W), 1) + W
    return(G)
  }
  Q <- function() {
    param <- interpret.theta()
    Q <- param$PREC *
      (Matrix::Diagonal(nrow(W), apply(W, 1, sum)) - param$alpha * W)
  }
}
```

```

    return(Q)
  }
  mu <- function() {
    return(numeric(0))
  }
  log.norm.const <- function() {
    val <- numeric(0)
    return(val)
  }
  log.prior <- function() {
    param <- interpret.theta()
    val <- -theta[1L] - 2 * log(1 + exp(-theta[1L]))
    ## PC prior
    val <- val + log(lambda/2) - theta[2L]/2 - (lambda * exp(-theta[2L]/2))
    ## Gamma(1, 5e-5), default prior:
    #val <- val + dgamma(exp(theta[2L]), shape = 1, rate = 5e-5, log = T) + theta[2L]
    ## Uniform prior on the standard deviation
    #val <- val - sum(theta[2L])/2 - k * log(2)
    return(val)
  }
  initial <- function() {
    return(c(0, 4))
  }
  quit <- function() {
    return(invisible())
  }
  if (as.integer(R.version$major) > 3) {
    if (!length(theta))
      theta = initial()
  }
  else {
    if (is.null(theta)) {
      theta <- initial()
    }
  }
  val <- do.call(match.arg(cmd), args = list())
  return(val)
}

PCAR.model <- function(...) INLA::inla.rgeneric.define(inla.rgeneric.PCAR.model, ...)

```

## Bibliography

- Gelman, Andrew, Jessica Hwang, and Aki Vehtari. 2014. “Understanding Predictive Information Criteria for Bayesian Models.” *Statistics and Computing* 24 (6): 997–1016. <https://doi.org/10.1007/S11222-013-9416-2>.
- Leroux, Brian G., Xingye Lei, and Norman Breslow. 2000. “Estimation of Disease Rates in Small Areas: A New Mixed Model for Spatial Dependence.” In *Statistical Models in Epidemiology, the Environment, and Clinical Trials*, edited by M. Elizabeth Halloran and Donald Berry, 179–91. New York, NY: Springer New York. [https://doi.org/10.1007/978-1-4612-1284-3\\_4](https://doi.org/10.1007/978-1-4612-1284-3_4).
- Palmí-Perales, Francisco, Virgilio Gómez-Rubio, and Miguel A. Martínez-Beneito. 2021. “Bayesian Multivariate Spatial Models for Lattice Data with INLA.” *Journal of Statistical Software* 98 (2): 1–29. <https://doi.org/10.18637/jss.v098.i02>.

Riebler, Andrea, Sigrunn H Sørbye, Daniel Simpson, and Håvard Rue. 2016. “An Intuitive Bayesian Spatial Model for Disease Mapping That Accounts for Scaling.” *Statistical Methods in Medical Research* 25 (4): 1145–65.