

# Explorative 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.

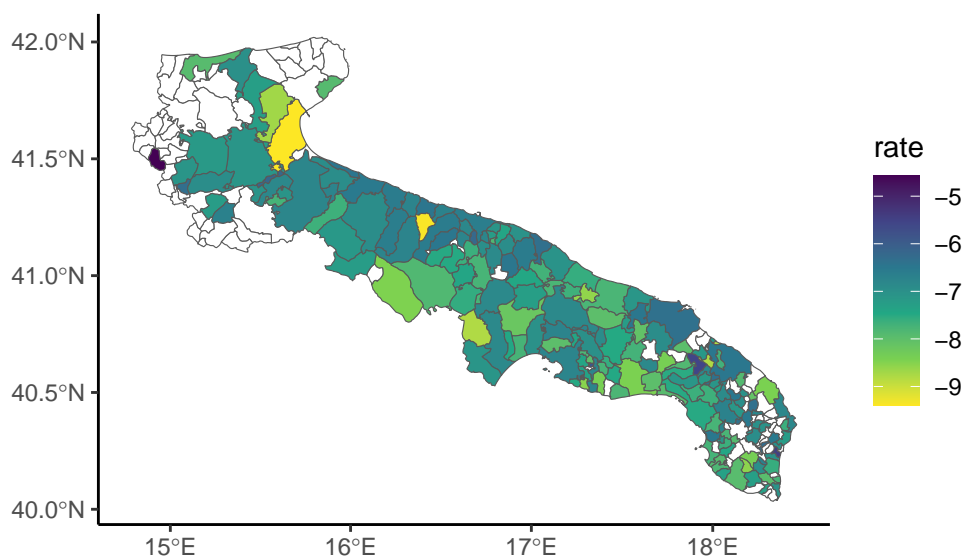
Here, we only take into account the violence reports which support centers actually take charge of, at the risk of underestimating the counts of gender-based violence cases. This choice is driven by the need of avoiding duplicated records, since e.g. it may happen that a support center redirects a victim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

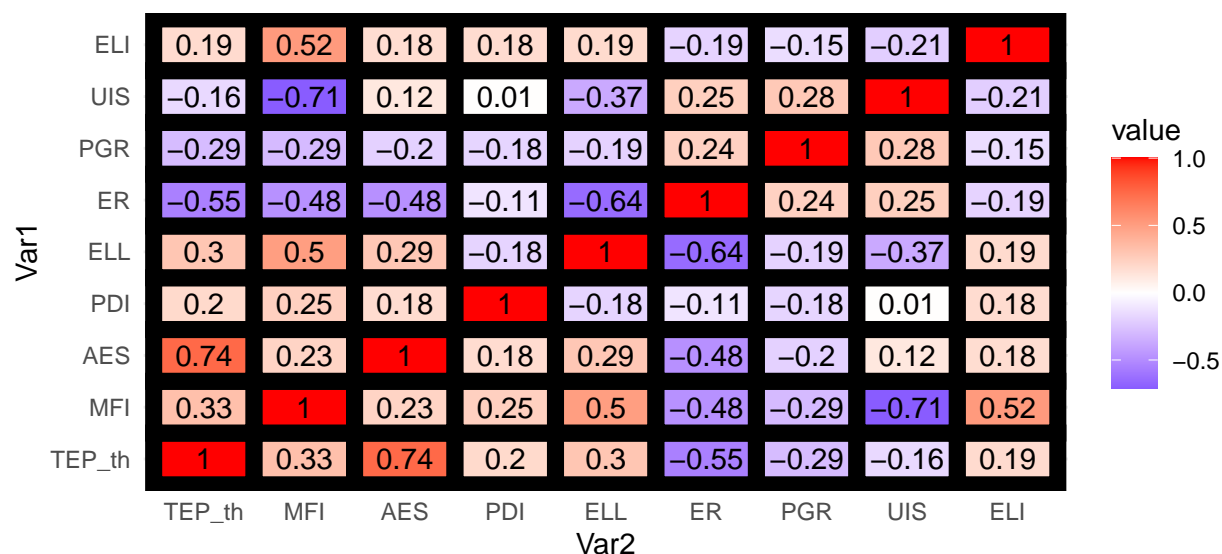
We model  $y$  via a simple Poisson GLM.

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.

- $d$ , i.e. the distance of each municipality from the closest municipality hosting a support center. Distance is measured by road travel time in minutes. For instance, the support center designated for the municipality of Adelfia (province of Bari, 3rd municipality in the dataset) is located in Triggiano (BA). Then,  $d_3$  denotes the travel time between Adelfia and Triggiano (17 minutes). In R,  $d$  is encoded as TEP\_th.
- 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:



Then, we implement a very simple forward selection algorithm. At each iteration, we add to the model the covariate allowing for the lowest BIC, until adding an additional covariate does not allow to reduce it anymore:

```
covariates <- colnames(X)[-1]
covs.in <- c()
BIC.min <- c()
```

```

while(length(covs.in) < length(covariates)){
  covs.out <- covariates[which(!covariates %in% covs.in)]

  BICs <- c()
  for(j in c(1:length(covs.out))) {
    formula.tmp <- paste0("N_ACC ~ 1 + offset(log(nn)) +",
                          paste(covs.in, collapse = "+"),
                          "+", covs.out[j])
    mod.tmp <- glm(formula.tmp, data = dd_con, family = "poisson")
    BICs[j] <- stats::BIC(mod.tmp)
  }
  BIC.min <- c(BIC.min, min(BICs))
  #if(length(BIC.min)>1 && BIC.min[length(BIC.min)] >= BIC.min[length(BIC.min)-1]){
  # break
  #} else{
    covs.in <- c(covs.in, covs.out[which.min(BICs)])
  #}
}

```

The optimal number of covariates (covariates in the model with minimum BIC) is two:

```
covs.in[c(1:which.min(BIC.min))]
```

```
## [1] "TEP_th" "AES"
```

However, a model with up to 4 covariates would have all significant regression coefficients:

```

summary(glm(
  as.formula(paste0("N_ACC ~ 1 +", paste(covs.in[c(1:5)], collapse = " + "))),
  family = "poisson", offset = log(nn), data = dd_con))$coefficients

```

##		Estimate	Std. Error	z value	Pr(> z )
##	(Intercept)	-7.48888909	0.04397971	-170.280546	0.000000e+00
##	TEP_th	-0.42037978	0.05285388	-7.953621	1.811370e-15
##	AES	-0.09567982	0.03408251	-2.807299	4.995878e-03
##	ELL	-0.06450769	0.03099637	-2.081137	3.742139e-02
##	UIS	-0.08895327	0.03397183	-2.618442	8.833222e-03
##	PGR	0.08186448	0.04354281	1.880092	6.009561e-02

In this case, the incidence of people with low educational level has a negative association with the accesses to support centers; the easiest interpretation would be that raising educational level would encourage reporting gender violence. Interpreting the negative association with the density of productive units becomes harder - and do not us forget that association does not imply a causal relationship. Lastly, the association with population growth rate (fifth covariate) appears positive, but the significance of this association has weaker evidence (p-value  $\sim 0.06$ ).

With more than fifth covariates, no additional valuable association can be found.

As we see, using the BIC as selection criterion, the most relevant covariates appear the distance indicators: distance from the closest support center and distance. These two variables have a high correlation (0.74).

In the remainder of this work, we will only focus on the two covariates  $d$  and AES

## Nonspatial regression

We regress the counts of accesses  $y$  to support centers on the distance from the former. Both covariates are scaled to zero mean an unit variance.

$$\frac{y_i}{P_i} \mid \eta_i \sim \text{Poisson}(e^{\eta_i}) \quad \text{where} \quad \eta_i = \beta_0 + \beta_d d_i + \beta_{AES} AES_i \quad (1)$$

Where  $P_i$  is the female population aged  $\geq 15$  in municipality  $i$ .

```
##
## Call:
## glm(formula = N_ACC ~ 1 + TEP_th + AES, family = "poisson", data = dd_con,
##      offset = log(nn))
##
## Coefficients:
##              Estimate Std. Error  z value Pr(>|z|)
## (Intercept) -7.49029      0.04374 -171.232  < 2e-16 ***
## TEP_th       -0.38633      0.04969  -7.775 7.55e-15 ***
## AES          -0.13625      0.03193  -4.267 1.98e-05 ***
## ---
## 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: 485.64  on 253  degrees of freedom
## AIC: 1058.2
##
## Number of Fisher Scoring iterations: 5
```

Both the distance from the closest support center and the distance from the closest infrastructural pole have a negative association with the frequency of accesses to support centers.

How do we interpret the regression coefficients? Keeping in mind we are working on the logarithm of the access rate, the standard deviations of explanatory variables, both expressed in minutes, are the following:

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

```
## [1] 15.74582
```

```
# Distance from closest infrastructural pole
attr(scale(Indicators$AES), "scaled:scale")
```

```
## [1] 15.47489
```

Hence e.g. each 15'45'' of distance of the a given municipality from the closest support center are associated with a decrease of 0.386 in the log-frequency at which women from that municipality access to support centers.

## Spatial regression

We plot the log-residuals  $\varepsilon$  of the regression model in equation 1, defined as  $\varepsilon := \ln y_i - \ln P_i - \ln \hat{y}_i$  being  $\hat{y}_i$  the fitted value.

```
## 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

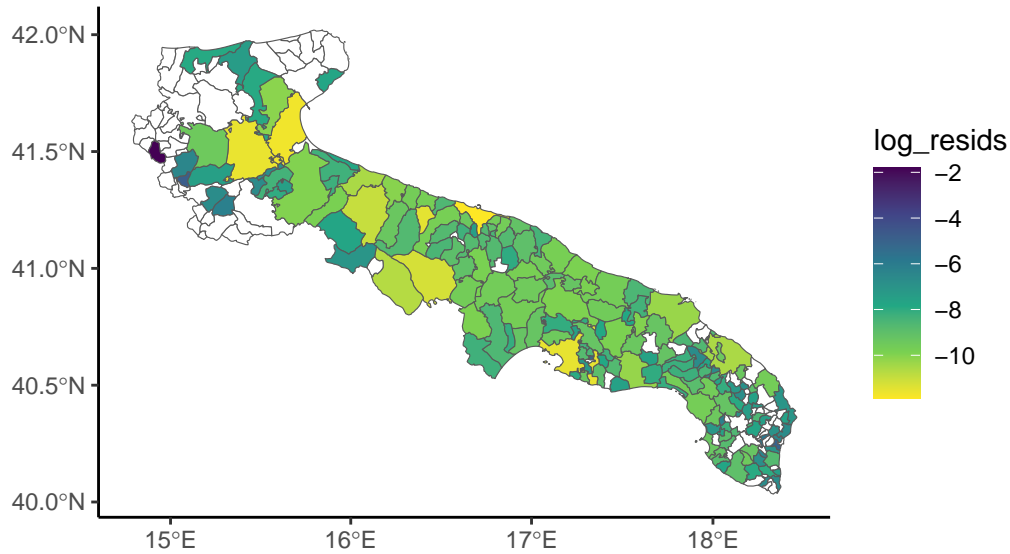


Figure 1: Log-residuals of glm regression using theoretical distance as explanatory variable

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 = 5.6902, p-value = 6.344e-09
## alternative hypothesis: greater
## sample estimates:
## Moran I statistic      Expectation      Variance
##      0.302415738      -0.005813953      0.002934223

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 = 4.0358, p-value = 2.721e-05
## alternative hypothesis: Expectation greater than statistic
## sample estimates:
## Geary C statistic      Expectation      Variance
##      0.751958417      1.000000000      0.003777415
```

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)

**ICAR model** Based on the autocorrelation evidence, though it has only been assessed for a subset of all municipalities, we try implementing a simple spatial model, say the intrinsic autoregressive:

$$\eta_i = \beta_0 + \beta_d d_i + \beta_{AES} AES_i + z_i \quad (2)$$

Where

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

Here  $w_{ij}$  is a generic element of the neighbourhood matrix of Apulia's municipalities, tanking value 1 if two municipalities  $i$  and  $j$  are neighbours and zero if either  $i$  and  $j$  are not neighbours or  $i = j$ ;  $d_i$  is the number of neighbours of municipality  $i$ , and  $\sigma^2$  is the scale parameter of the spatial effect.  $n = 256$  and the remainder of the model follows the same notation as eq. 1.

On a first stance, we will let the Laplace Approximation approximate posterior distributions of the quantities of interest (instead of e.g. relying on MCMC)

We set a PC-prior on the random effect to avoid overfitting. The `f()` argument sets the prior model on the spatial latent effect. Arguments `num.threads` and `internal.opt` are used for reproducibility. We use a PC-prior with rate parameter  $\lambda = 1.5$ , such that  $\text{Prob}(\sigma > \lambda) = 0.01$

```
cav_icar_INLA <- inla(
  N_ACC ~ 1 + TEP_th + AES + f(ID, model = "besag", graph = W_con,
    scale.model = T, prior = "pc.prec", param = c(1.5, 0.01)),
  family = "poisson", offset = log(nn), data = dd_con,
  num.threads = 1, control.compute = list(internal.opt = F, cpo = T, waic = T),
  inla.mode = "classic", control.inla = list(strategy = "laplace"),
  verbose = F) # better
```

We compare the spatial and the nonspatial model (for coparability it has been re-run with R-INLA) through the Watanabe-Akaike information criterion

```
cav_nosp_INLA$waic$waic # Higher WAIC
```

```
## [1] 1064.348
```

```
cav_nosp_INLA$waic$p.eff # Low complexity
```

```
## [1] 8.605647
```

```
cav_icar_INLA$waic$waic # Lower WAIC
```

```
## [1] 958.8593
```

```
cav_icar_INLA$waic$p.eff # High complexity
```

```
## [1] 72.69854
```

Which suggests us the spatial model is preferable; model complexity increases (`p.eff` denotes the number of free/unconstrained parameters) but is still outweighed by better fitting. Here we show some summaries for the spatial model:

```
cav_icar_INLA$summary.fixed
```

```
##               mean          sd 0.025quant  0.5quant  0.975quant      mode
## (Intercept) -7.6230980 0.06223487 -7.7493003 -7.6217593 -7.50479148 -7.6216469
## TEP_th      -0.3311004 0.08771165 -0.5044677 -0.3307689 -0.15964244 -0.3307749
## AES         -0.1760451 0.07353991 -0.3223210 -0.1754928 -0.03291102 -0.1754966
```

```
##                                kld
## (Intercept) 2.061461e-05
## TEP_th      5.250088e-08
## AES         5.990678e-08
```

We still see distance from the closest support center has a negative effect on the accesses. A tentative interpretation we may draw from this result is that the distance from support centers may tend to discourage women from reporting gender-based violence, and the occurrence of gender-based violence may be underestimated in areas distant from support centers.

The spatial effect is to be interpreted as the spatial variation not explained by model components. Here we show the posterior summaries of the marginal variance parameter ( $\sigma^2$ ) of the spatial effect

```
inla.zmarginal(inla.tmargin(
  fun = function(X) 1/X,
  marginal = cav_icar_INLA$marginals.hyperpar[[1]]))
```

```
## Mean          0.657288
## Stdev          0.182537
## Quantile 0.025 0.350284
## Quantile 0.25  0.526114
## Quantile 0.5   0.639436
## Quantile 0.75  0.76874
## Quantile 0.975 1.06491
```

**PCAR model** Now, 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 assume a 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  $\alpha$ :

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

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, ...)

```

Then the model is run

```

cav_pcar_INLA <- inla(N_ACC ~ 1 + TEP_th + AES +
  f(ID, model = PCAR.model(W = W_con, k = 1, lambda = 1.5)),
  family = "poisson", offset = log(nn), data = dd_con,
  num.threads = 1, control.compute =
    list(internal.opt = F, cpo = T, waic = T),
  inla.mode = "classic", control.inla = list(strategy = "laplace"),
  control.predictor = list(compute = T),
  verbose = T)

```



We compare the WAICs of the two spatial models

```
cav_icar_INLA$waic$waic # Higher WAIC
```

```
## [1] 958.8593
```

```
cav_icar_INLA$waic$p.eff
```

```
## [1] 72.69854
```

```
cav_pcar_INLA$waic$waic # Lower WAIC
```

```
## [1] 940.7756
```

```
cav_pcar_INLA$waic$p.eff # Complexity is not even raised
```

```
## [1] 71.39892
```

The lower WAIC suggests us that the PCAR, which does not even appear to be more complex although the additional parameter  $\alpha$ . We show the posterior summaries of the regression coefficients

```
cav_pcar_INLA$summary.fixed
```

```
##              mean          sd 0.025quant  0.5quant  0.975quant      mode
## (Intercept) -7.6187902 0.07249439 -7.7656954 -7.6173246 -7.48032377 -7.6172846
## TEP_th      -0.3634936 0.08235752 -0.5261945 -0.3631800 -0.20256749 -0.3631796
## AES         -0.1978873 0.06704965 -0.3307792 -0.1975221 -0.06703311 -0.1975143
##              kld
## (Intercept) 4.183002e-06
## TEP_th      5.357745e-08
## AES         5.227461e-08
```

Spatial autocorrelation does not appear particularly strong. Additionally, the credible interval appears uncannily wide:

```
## Mean          0.470016
## Stdev         0.222425
## Quantile 0.025 0.0809099
## Quantile 0.25  0.290043
## Quantile 0.5   0.47005
## Quantile 0.75  0.646321
## Quantile 0.975 0.869732
```

Lastly, we plot the estimated PCAR field:

Spatial pattern does not appear particularly strong. We may need for a model accounting for random noise instead of this one.

**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  $\in [0, 1]$ .

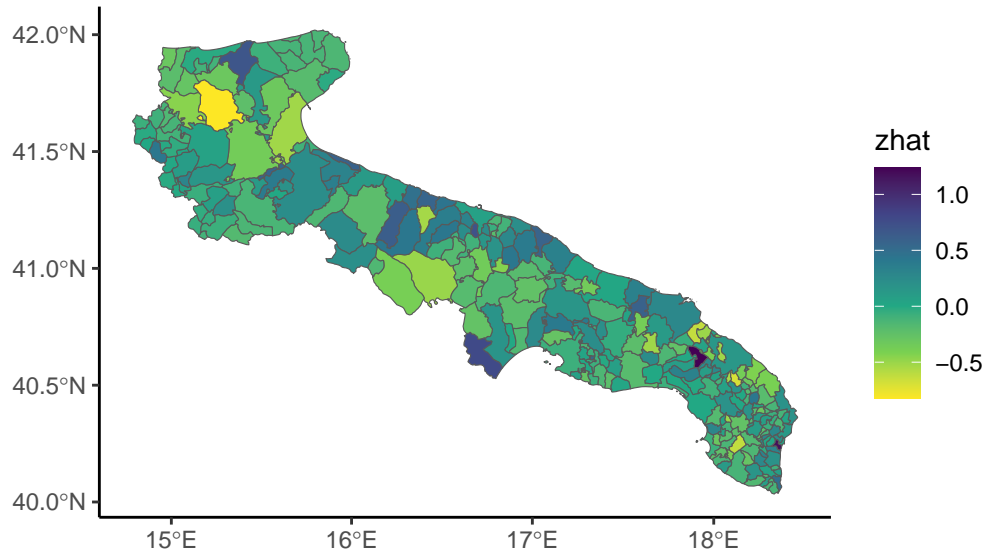


Figure 2: Posterior expectation of the spatial latent effect, PCAR model

```
cav_bym_INLA <- inla(
  N_ACC ~ 1 + TEP_th + AES + f(ID, model = "bym2", graph = W_con,
    scale.model = T, prior = "pc.prec", param = c(1.5, 0.01)),
  family = "poisson", offset = log(nn), data = dd_con,
  num.threads = 1, control.compute = list(internal.opt = F, cpo = T, waic = T),
  inla.mode = "classic", control.inla = list(strategy = "laplace"),
  verbose = F)
```

First, we take a look at the WAIC of the BYM model

```
cav_bym_INLA$waic$waic
```

```
## [1] 940.1248
```

```
cav_bym_INLA$waic$p.eff
```

```
## [1] 70.61776
```

As we can see, this model has a lower WAIC than the PCAR. This reduction is due to the smaller number of free parameters (which actually decreases more than the overall WAIC, implying a slight and more than outweighed loss in fitting quality).

Next we show the hyperparameters posterior

```
cav_bym_INLA$summary.hyperpar
```

```
##               mean          sd 0.025quant  0.5quant 0.975quant
## Precision for ID 4.6478440 1.2605358 2.737540491 4.46060350 7.6596208
## Phi for ID       0.1335076 0.1427875 0.003234544 0.08001105 0.5348934
##               mode
## Precision for ID 4.079037152
## Phi for ID       0.005174409
```

As we see, the mixing parameter is low, and its estimation is affected by strong uncertainty. Based on this distribution, we can indeed argue that most of the variation in  $y$  not explained by covariates has to be attributed to noise.

Finally, we show posterior summaries for covariates effects:

```
cav_bym_INLA$summary.fixed
```

```
##              mean          sd 0.025quant  0.5quant  0.975quant      mode
## (Intercept) -7.6160297 0.06471917 -7.7468645 -7.6147182 -7.49279285 -7.6146433
## TEP_th      -0.3411954 0.08641847 -0.5104081 -0.3414229 -0.17063290 -0.3414100
## AES         -0.1963900 0.07040219 -0.3356727 -0.1960946 -0.05876719 -0.1960922
##              kld
## (Intercept) 8.234444e-06
## TEP_th      4.855938e-08
## AES         4.346022e-08
```

Interpretation does not differ too much from the GLM estimates: both variables have a negative impact on the access rate to support centers. Estimated effects are lower in absolute values.

## 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.

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.

## Bibliography

- 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.