



## **DClusterm: Model-based detection of disease clusters**

**Virgilio Gómez-Rubio**

Universidad de Castilla-La Mancha

**Paula Moraga-Serrano**

London School of Hygiene  
& Tropical Medicine

**John Molitor**

Oregon State University

**Barry Rowlingson**

Lancaster University

---

### **Abstract**

The detection of regions with unusual high risk plays an important role in disease mapping and the analysis of Public Health data. In particular, the detection of groups of areas (i.e., clusters) where the risk is significantly high is often conducted by Public Health authorities.

Many methods have been proposed for the detection of disease clusters, most of them based on moving windows, such as, Kulldorff's Spatial Scan Statistics (SSS). Here we describe a model-based approach for the detection of disease clusters implemented in the **DClusterm** package. Our model-based approach is based on representing a large number of possible clusters by dummy variables and then fitting many generalized linear models to the data where these covariates are included one at a time. Cluster detection is done by performing a variable or model selection among all fitted models using different criteria.

Because of our model-based approach, cluster detection can be performed using different types of likelihoods and latent effects. We cover the detection of spatial and spatio-temporal clusters, as well as how to account for covariates, deal with zero-inflated datasets and overdispersion in the data.

*Keywords:* disease cluster, spatial statistics, R.

---

## **1. Introduction**

The analysis of epidemiological data at small area level often involves accounting for possible risk factors and other important covariates using different types of regression models. How-

ever, it is not uncommon that after a number of covariates have been accounted for, residuals show a spatial distribution that defines some groups of areas with unusual high epidemiological risk. Hence, in many occasions it is not clear whether all spatial risk factors have been included in our model.

Public health data are often aggregated over small administrative areas because of confidentiality issues, but it is not uncommon that individual data are available. Generalised Linear Models (GLM, ) are a common framework for disease mapping to model aggregated and individual data. GLMs not only model Poisson or Binomial responses, but they can also link the outcome to a linear predictor on the covariates (and, possibly, other effects). However, until recently, it was not clear how to use GLMs to detect clusters of disease, i.e., a group of contiguous areas with significant high risk.

In order to detect disease clusters, probably the most widely used method is the one proposed by [Kulldorff \(1997\)](#). This is called the Spatial Scan Statistic and it will find the most likely cluster. Significance is assessed via a Monte Carlo test using a test statistic based on a likelihood ratio test for the following hypotheses:

$$\begin{aligned} H_0 : \theta_z &= \theta_{\bar{z}} \\ H_1 : \theta_z &> \theta_{\bar{z}} \end{aligned}$$

Here,  $z$  represents a cluster (i.e., a set of contiguous areas),  $\theta_z$  the relative risk in the cluster and  $\theta_{\bar{z}}$  the relative risk outside the cluster. Many different clusters are tested in turn. The most likely cluster is the one with the highest value of the test statistic. Then a Monte Carlo test is used to compute the p-value of the most likely cluster.

In this paper we will summarise the work by several authors that have established a link between GLMs and SSS, so that the detection of disease clusters is approached from a regression point of view. As described later, this will involve fitting many different GLMs for which dummy variables that represent possible clusters are included one at a time. Cluster detection is based on selecting a number of dummy cluster variables using variable selection methods. Furthermore, we will describe how these methods have been implemented in the **DClusterm** package for the R software.

This paper is organised as follows. Section 2 will introduce the link between GLM and SSS. Next, in Section 3 we describe how to extend these ideas to detect clusters in space and time. The detection of disease clusters for zero-inflated data is discussed in Section 4. Section 5 shows how to include random effects in the detection of disease clusters. A multivariate approach for the detection of disease clusters of two diseases has been included in Section 6. Finally, a discussion and some final remarks are provided in Section 7.

## 2. Generalised Linear Models for cluster detection

{sec:GLM}

[Jung \(2009\)](#); [Zhang and Lin \(2009\)](#) provide an explicit link between GLMs and the SSS, and show that the test statistic for a given cluster is equivalent to fitting a Generalised Linear Model using a cluster variable as a predictor. This cluster variable is a dummy variable which is 1 for the areas in the cluster and 0 for the areas outside the cluster.

Firstly, given that we are using GLM's we could include covariates in the model. For example, for a Poisson model with expected counts  $E_i$  we could have:

$$O_i \sim Po(E_i\theta_i)$$

$$\log(\theta_i) = \log(E_i) + \alpha + \beta x_i$$

Fitting this model will provide estimates  $\hat{\alpha}$  and  $\hat{\beta}$ . This will account for the (spatial) effects of the covariates. In order to include the cluster variable the effects of the covariates will be kept fixed. Hence, the clusters covariates will be used in a model with fixed coefficients for the covariates:

$$\log(\theta_i) = \log(E_i) + \hat{\alpha} + \hat{\beta}x_i + \gamma CLUSTER_i$$

This means that the offset now is  $\log(E_i) + \hat{\alpha} + \hat{\beta}x_i$ .  $\gamma$  is a measure of the difference of the risk in the cluster. We are only interested in clusters whose coefficient is higher than 0 (i.e., increased risk), hence those with a significant negative coefficient will be ignored.

Testing different clusters will produce many different cluster covariates. We can use model selection techniques to select the most important cluster in the area. In particular, the log-likelihood can be used to compare the model with the cluster variable to the null model (i.e., the one with the covariates only). Note that we are interested in clusters with a high risk and, because of that, we are only interested in clusters whose associated coefficient is significantly higher than zero.

Regarding the effect of the covariates, it is possible to perform a cluster detection without considering covariates in the model. Then a cluster detection accounting for the covariates will likely provide a different number of clusters. By comparing the clusters detected in both cases we will be able to find what clusters are linked to underlying risk factors included in the model and what clusters remain unexplained by the covariates. In the examples that we included in this paper we will always consider both scenarios to better understand how cluster detection works with these methods.

Bilancia and Demarinis (2014); Gómez-Rubio, Moraga, and Molitor (2015) describe a similar approach to the detection of disease cluster using Bayesian hierarchical models. The Integrated Nested Laplace Approximation is used in both cases for model fitting as it provides computational benefits over other computationally expensive methods, such as Markov Chain Monte Carlo.

## 2.1. Leukemia in upstate New York

The NY8 dataset is available in package `DCluster` and it provides cases of leukemia in different census tracts in upstate New York. This data set has been analysed by several authors (Waller, Turnbull, Clark, and Nasca 1992; Waller and Gotway 2004). The location of leukemia is thought to be linked to the use of Trichloroethene (TCE) by several companies in the area. Figure 1 shows the Standardised Mortality Ratios of the census tracts and the locations of the industries using TCE.

In order to measure exposure, the inverse of the distance to the nearest TCE site has been used (PEXPOSURE). In addition, two other socioeconomic covariates have been used: the percentage of people aged 65 or more (PCTAGE65P) and the percentage of people who own their home (PCTOWNHOME).

This dataset is included in package **DClusterm** as `NY8`. Hence, our first action is to load some required packages and the dataset itself.

```
> library(DClusterm)
> library(snowfall)
> library(xts)
> data(NY8)
```

A number of cases could not be linked to their actual location and they were distributed uniformly over the study area, making the counts real numbers instead of integers. We have rounded these values as we intend to use a Poisson likelihood for the analysis. Furthermore, expected counts are computed using the overall incidence ratio (total number of cases divided by the total population). Age-sex standardisation is not possible in this case as this information is not available in our dataset.

```
> NY8$Observed <- round(NY8$Cases)
> NY8$Expected <- NY8$POP8 * sum(NY8$Observed)/sum(NY8$POP8)
> NY8$SMR <- NY8$Observed/NY8$Expected
> NY8$x <- coordinates(NY8)[, 1]
> NY8$y <- coordinates(NY8)[, 2]
```

Finally, a `STFDF` object is created to store all the data. Functions in **DClusterm** will take object for space-time data as defined in package **spacetime**. Note that in this case we do not have a truly space-time dataset.

```
> NY8st <- STFDF(as(NY8, "SpatialPolygons"), xts(1, as.Date("1972-01-01")),
+   NY8@data, endTime = as.POSIXct(strptime(c("1972-01-01"),
+   "%Y-%m-%d"), tz = "GMT"))
```

## 2.2. Cluster detection

### *Cluster detection with no covariates*

First of all, a model with no covariates will be fitted and used as a baseline for model fitting. For example, other models can be compared to this one (for example, using the AIC or the log-likelihood) to assess whether they provide a better fit.

```
> m0 <- glm(Observed ~ offset(log(Expected)) + 1, family = "poisson",
+   data = NY8)
```

Cluster detection will use the previous model and new cluster dummy variables will be included, one at a time, to test for a large number of clusters.

Function `DetectClustersModel()` will take the baseline model (using argument `model0`), create the cluster dummy variables and test them in turn. Then, those clusters with a highest significance will be reported.

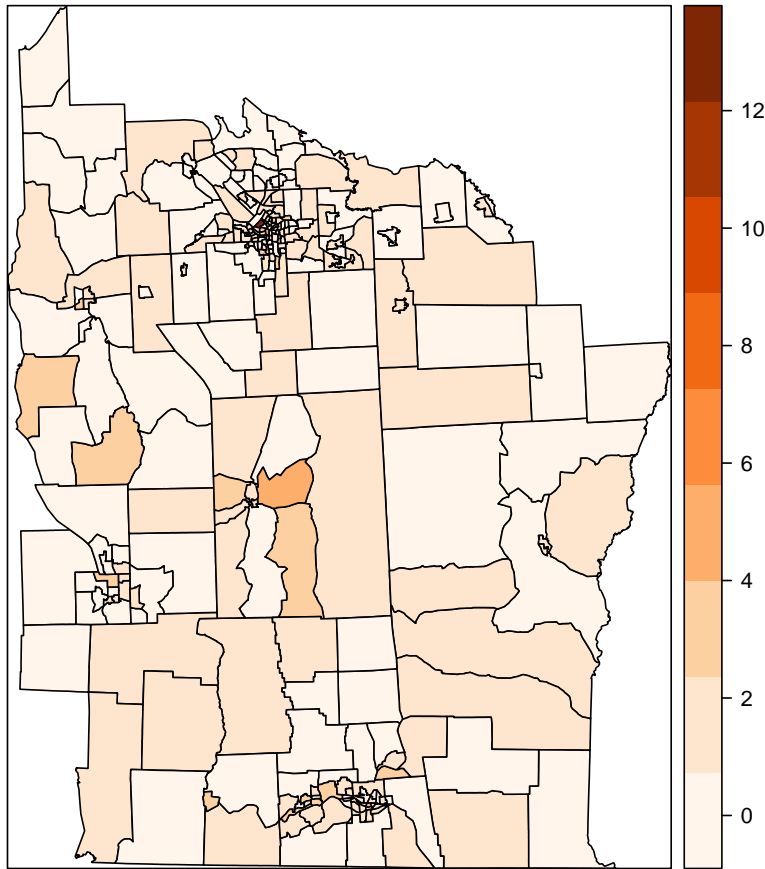


Figure 1: **\*\*INCLUDE TCE LOCATIONS AND THREE OTHER COVARIATES\*\*** SMR of the incidence of Leukemia in upstate New York.

{fig:NYmap}

Argument `thegrid` will take a 2-column `data.frame` (with names `x` and `y`) with the centres of possible clusters. If the grid of cluster centres is not defined, then a rectangular grid is used with a distance between adjacent points defined by argument `step`. Dummy cluster variables are created around these points by adding areas to the cluster until a certain percentage of the population has been reached (defined by argument `fractpop`) or until a certain distance about the centre (defined by argument `radius`). When testing for significant cluster variables, argument `alpha` defines the significance level.

`DetectClustersModel()` can detect spatial and spatio-temporal clusters, that is why its first argument is a space-time object. The type of clusters that are investigated is defined by argument `typeCluster`. In the example we have used `typeCluster = "S"`.

Other options include the number of CPUs to be used to test for clusters in parallel (argument `numCPUS`) and the number of replicates for Monte Carlo tests (argument `R`) if cluster assessment is done by simulation. By default, Monte Carlo tests are not used.

In the following example, to reduce the computational burden, we have only looked for clusters around 5 areas (whose rows in `NY8` are defined in variable `idxcl`). In a real application we

advice the use of all locations (area centroids or actual locations of individual data).

```
> idxcl <- c(120, 12, 89, 139, 146)
> cl0 <- DetectClustersModel(NY8st, thegrid = as.data.frame(NY8)[idxcl,
+   c("x", "y")], fractpop = 0.15, alpha = 0.05, radius = Inf,
+   step = NULL, typeCluster = "S", R = NULL, numCPUS = 4, model0 = m0)
```

Below is a summary of the clusters detected. Dates can be ignored as this is a purely spatial cluster. In the case of spatio-temporal clusters, the dates shown define the temporal range of the cluster. Values **x** and **y** defined the cluster centre, **size** is the number of areas (or individuals) in the cluster, **statistic** represents the point estimate of the associated cluster coefficient. Also, note that only clusters with a lower **pvalue** than argument **alpha** are returned. **cluster** indicates whether the cluster is a significant one. Finally, note how detected cluster are order by increasing value of **pvalue**, so that most significant clusters are reported first.

```
> cl0
```

|     | x        | y       | size | minDateCluster      | maxDateCluster      | statistic |
|-----|----------|---------|------|---------------------|---------------------|-----------|
| 11  | 424728.9 | 4661404 | 39   | 1972-01-01 01:00:00 | 1972-01-01 01:00:00 | 8.044846  |
| 88  | 409430.4 | 4720092 | 9    | 1972-01-01 01:00:00 | 1972-01-01 01:00:00 | 6.967107  |
| 119 | 404710.7 | 4768346 | 24   | 1972-01-01 01:00:00 | 1972-01-01 01:00:00 | 3.254824  |

|     | cluster | pvalue       |
|-----|---------|--------------|
| 11  | TRUE    | 0.0000604120 |
| 88  | TRUE    | 0.0001893208 |
| 119 | TRUE    | 0.0107290781 |

The centre of the clusters detected are shown in Figure 2. Because of the lack of adjustment for covariates these clusters show regions of high risk based on the raw data (observed and expected counts) alone.

### *Cluster detection after adjusting for covariates*

Similarly, clusters can be detected after adjusting for significant risk factors. First of all, we will fit a Poisson regression with the 3 covariates mentioned earlier. As it can be seen, all three are significant:

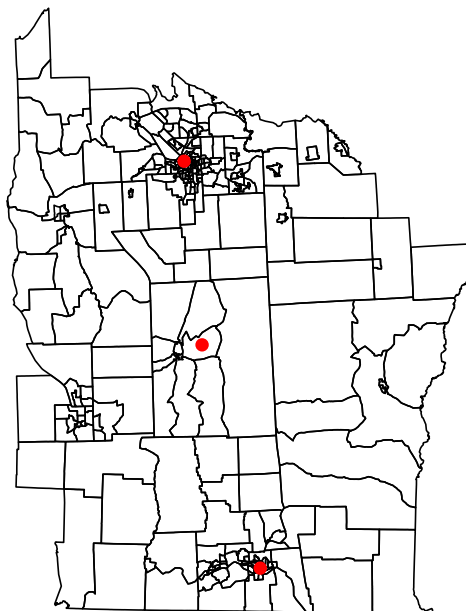
```
> m1 <- glm(Observed ~ offset(log(Expected)) + PCTOWNHOME + PCTAGE65P +
+   PEXPOSURE, family = "poisson", data = NY8)
> summary(m1)
```

Call:

```
glm(formula = Observed ~ offset(log(Expected)) + PCTOWNHOME +
    PCTAGE65P + PEXPOSURE, family = "poisson", data = NY8)
```

Deviance Residuals:

| Min | 1Q | Median | 3Q | Max |
|-----|----|--------|----|-----|
|-----|----|--------|----|-----|



{fig:NYc10}

Figure 2: Clusters detected when no covariates are included in the model.

-2.9099   -1.1294   -0.1768   0.6385   3.2426

Coefficients:

|             | Estimate | Std. Error | z value | Pr(> z ) |     |
|-------------|----------|------------|---------|----------|-----|
| (Intercept) | -0.65507 | 0.18550    | -3.531  | 0.000413 | *** |
| PCTOWNHOME  | -0.36472 | 0.19316    | -1.888  | 0.058998 | .   |
| PCTAGE65P   | 4.05031  | 0.60559    | 6.688   | 2.26e-11 | *** |
| PEXPOSURE   | 0.15141  | 0.03165    | 4.784   | 1.72e-06 | *** |

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for poisson family taken to be 1)

Null deviance: 459.05 on 280 degrees of freedom  
 Residual deviance: 384.01 on 277 degrees of freedom  
 AIC: 958.97

Number of Fisher Scoring iterations: 5

As the three covariates are significant, the expected number of cases will be different now and the detected clusters may be different in this case. Cluster detection is performed as in the previous example, but now we use the model that adjusts for covariates instead:

```
> cl1 <- DetectClustersModel(NY8st, thegrid = as.data.frame(NY8)[idxcl,
+   c("x", "y")], fractpop = 0.15, alpha = 0.05, typeCluster = "S",
+   R = NULL, numCPUS = 4, model0 = m1)

> cl1
```

|     | x        | y       | size | minDateCluster      | maxDateCluster      | statistic |
|-----|----------|---------|------|---------------------|---------------------|-----------|
| 88  | 409430.4 | 4720092 | 9    | 1972-01-01 01:00:00 | 1972-01-01 01:00:00 | 5.861204  |
| 119 | 404710.7 | 4768346 | 20   | 1972-01-01 01:00:00 | 1972-01-01 01:00:00 | 3.160591  |

|     | cluster | pvalue       |
|-----|---------|--------------|
| 88  | TRUE    | 0.0006175202 |
| 119 | TRUE    | 0.0119304026 |

Figure 3 shows the clusters detected after adjusting for covariates. Compared to the example with no covariate adjustment, one cluster has disappeared. Hence, that cluster has been explained by the effect of the covariates. Another cluster is a bit smaller in size, which means that covariate only explain a small part of it. The most significant cluster remains the same. In all cases, cluster significance has been reduced by the effect of the covariates.

### 3. Spatio-temporal clusters

Jung (2009) discusses how to extend model-based approaches for the detection of spatial disease clusters to space and time. Gómez-Rubio *et al.* (2015) propose the following model:

$$\log(\mu_{i,t}) = \log(E_{i,t}) + \gamma_j c_{i,t}^{(j)} \quad (1)$$

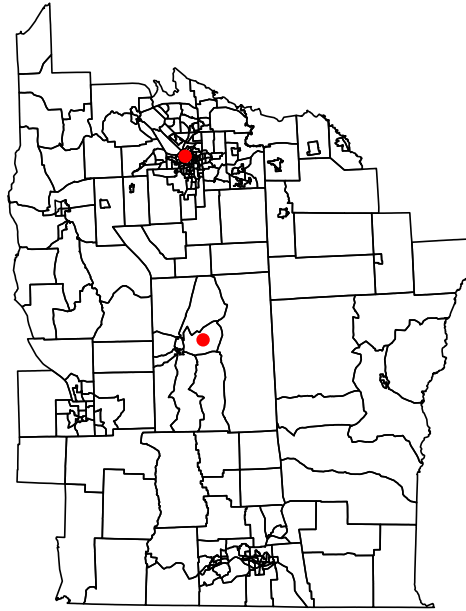
where  $\mu_{i,t}$  is the mean of area  $i$  at time  $t$  and  $c_{i,t}^{(j)}$  a cluster dummy variable for cluster  $j$ .

Note how now data are indexed according to space and time. Dummy cluster variables are defined as in the spatial case, by considering areas in the cluster according to their distance to the cluster centre, for data within a particular time period. When defining a temporal cluster, areas are aggregated using all possible temporal windows up to a predefined temporal range.

#### 3.1. Brain Cancer in New Mexico

The **brainNM** dataset (included in **DClusterm**) contains yearly cases of brain cancer in New Mexico from 1973 to 1991 (inclusive) in a **spacetime** object. The data set has been taken from the SatScan website and the area boundaries from the U.S. Census Bureau. In addition, the location of Los Alamos National Laboratory has been included (from the Wikipedia). Inverse distance to this site can be used to test for increased risk in the areas around the Laboratory as no other covariates are available.





{fig:NYc11}

Figure 3: Clusters detected after adjusting for covariates.

```
> data(brainNM)
```

Expected counts have been obtained using age and sex standardisation over the whole period of time. Hence, yearly differences are likely to be seen when plotting the data. Standardised Mortality Ratios have been plotted in Figure 4.

### 3.2. Cluster detection

#### *Cluster detection with no covariates*

Similarly as in the purely spatial case, a Poisson regression with no covariates will be fitted:

```
> m0 <- glm(Observed ~ offset(log(Expected)) + 1, family = "poisson",
+   data = brainst)
> summary(m0)
```

Call:

```
glm(formula = Observed ~ offset(log(Expected)) + 1, family = "poisson",
```

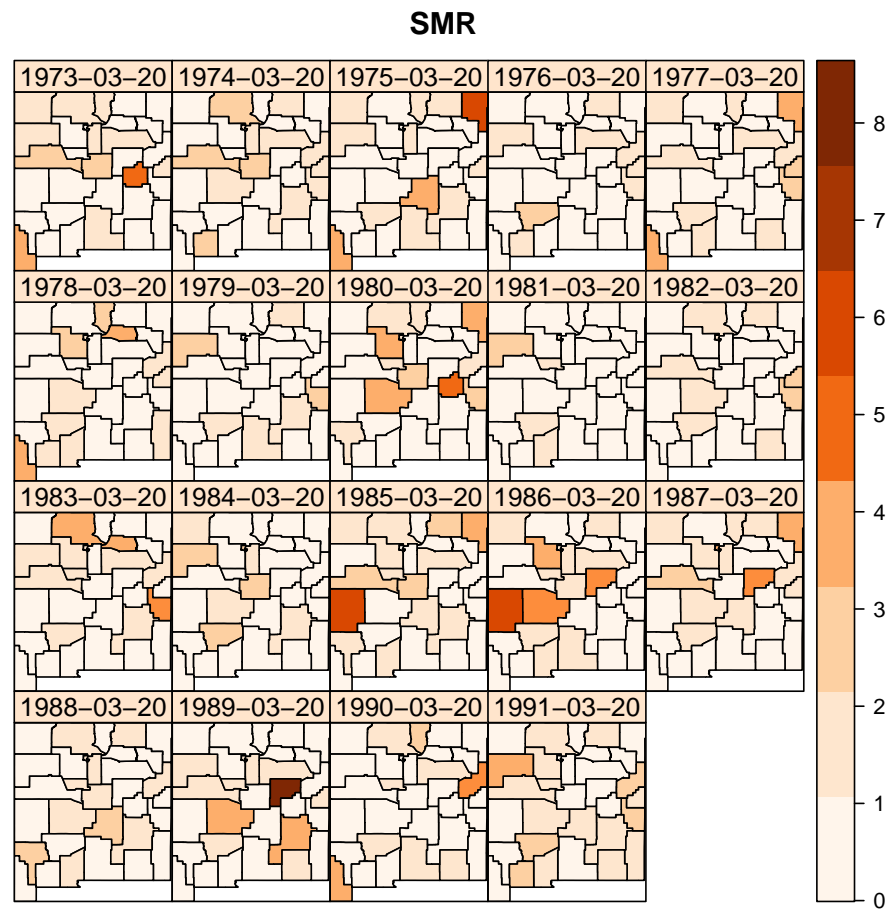


Figure 4: Standardised Mortality Ratios of brain cancer in New Mexico.

{fig:NMSMR}

```
data = brainst)

Deviance Residuals:
    Min       1Q   Median       3Q      Max
-2.4874  -0.9998  -0.4339   0.3773   3.1321

Coefficients:
              Estimate Std. Error z value Pr(>|z|)
(Intercept) 2.834e-16  2.917e-02      0      1

(Dispersion parameter for poisson family taken to be 1)

Null deviance: 631.64  on 607  degrees of freedom
Residual deviance: 631.64  on 607  degrees of freedom
AIC: 1585.6

Number of Fisher Scoring iterations: 5
```

Cluster detection with function `DetectClustersModel()` takes now arguments `minDateUser` and `maxDateUser` to define the minimum and maximum times that are considered when looking for clusters. `typeCluster = "ST"` is used to look for spatio-temporal clusters. **\*\*ANYTHING ELSE ABOUT HOW S-T CLUSTERS ARE DEFINED?\*\***

```
> c10 <- DetectClustersModel(brainst, coordinates(brainst@sp),
+   minDateUser = "1985-01-01", maxDateUser = "1989-01-01", fractpop = 0.15,
+   alpha = 0.05, typeCluster = "ST", R = NULL, numCPUS = 4,
+   model0 = m0)

> nrow(c10)

[1] 180

> c10[1:5, ]
```

|      | x         | y        | size | minDateCluster      | maxDateCluster      | statistic |
|------|-----------|----------|------|---------------------|---------------------|-----------|
| 0286 | -106.3073 | 35.86930 | 3    | 1986-03-20 01:00:00 | 1988-03-20 01:00:00 | 7.493492  |
| 0496 | -105.9761 | 35.50684 | 2    | 1986-03-20 01:00:00 | 1988-03-20 01:00:00 | 6.438221  |
| 0531 | -106.9303 | 34.00725 | 9    | 1985-03-20 01:00:00 | 1986-03-20 01:00:00 | 6.378992  |
| 0498 | -105.9761 | 35.50684 | 2    | 1987-03-20 01:00:00 | 1988-03-20 01:00:00 | 6.331113  |
| 0288 | -106.3073 | 35.86930 | 2    | 1987-03-20 01:00:00 | 1988-03-20 01:00:00 | 6.331113  |

|      | cluster | pvalue       |
|------|---------|--------------|
| 0286 | TRUE    | 0.0001082553 |
| 0496 | TRUE    | 0.0003327442 |
| 0531 | TRUE    | 0.0003544929 |
| 0498 | TRUE    | 0.0003731179 |
| 0288 | TRUE    | 0.0003731179 |

### *Cluster detection after adjusting for covariates*

We will use the inverse of the distance to Los Alamos National Laboratory as a covariate.

```
> dst <- spDistsN1(coordinates(brainst@sp), losalamos, TRUE)
> nyears <- length(unique(brainst@data$Year))
> brainst@data$IDLANL <- rep(1/dst, nyears)

> m1 <- glm(Observed ~ offset(log(Expected)) + IDLANL, family = "poisson",
+   data = brainst)
> summary(m1)
```

Call:

```
glm(formula = Observed ~ offset(log(Expected)) + IDLANL, family = "poisson",
    data = brainst)
```

Deviance Residuals:

| Min     | 1Q      | Median  | 3Q     | Max    |
|---------|---------|---------|--------|--------|
| -2.4832 | -0.9982 | -0.4280 | 0.3775 | 3.1424 |

Coefficients:

|             | Estimate  | Std. Error | z value | Pr(> z ) |
|-------------|-----------|------------|---------|----------|
| (Intercept) | -0.005721 | 0.029897   | -0.191  | 0.848    |
| IDLANL      | 0.338194  | 0.364900   | 0.927   | 0.354    |

(Dispersion parameter for poisson family taken to be 1)

Null deviance: 631.64 on 607 degrees of freedom  
 Residual deviance: 630.84 on 606 degrees of freedom  
 AIC: 1586.8

Number of Fisher Scoring iterations: 5

```
> cl1 <- DetectClustersModel(brainst, coordinates(brainst@sp),
+   fractpop = 0.15, alpha = 0.05, minDateUser = "1988-01-01",
+   maxDateUser = "1989-01-01", typeCluster = "ST", R = NULL,
+   numCPUS = 4, model0 = m1)
```

```
> nrow(cl1)
```

```
[1] 6
```

```
> cl1[1:5, ]
```

|     | x         | y          | size | minDateCluster      | maxDateCluster      | statistic |
|-----|-----------|------------|------|---------------------|---------------------|-----------|
| 049 | -105.9761 | 35.50684   | 2    | 1988-03-20 01:00:00 | 1988-03-20 01:00:00 | 2.433451  |
| 028 | -106.3073 | 35.86930   | 2    | 1988-03-20 01:00:00 | 1988-03-20 01:00:00 | 2.433451  |
| 057 | -105.8508 | 34.64048   | 2    | 1988-03-20 01:00:00 | 1988-03-20 01:00:00 | 2.431998  |
| 013 | -106.8328 | 32.35265   | 17   | 1988-03-20 01:00:00 | 1988-03-20 01:00:00 | 2.010047  |
| 027 | -105.4592 | 33.74524   | 3    | 1988-03-20 01:00:00 | 1988-03-20 01:00:00 | 2.007057  |
|     | cluster   | pvalue     |      |                     |                     |           |
| 049 | TRUE      | 0.02737662 |      |                     |                     |           |
| 028 | TRUE      | 0.02737662 |      |                     |                     |           |
| 057 | TRUE      | 0.02742274 |      |                     |                     |           |
| 013 | TRUE      | 0.04496121 |      |                     |                     |           |
| 027 | TRUE      | 0.04512090 |      |                     |                     |           |

We can easily display the most significant cluster as follows:

```
> stcl <- get.stclusters(brainst, cl0)
> brainst$CLUSTER <- 0
> brainst$CLUSTER[stcl[[1]]] <- 1
```

```
> print(stplot(brainst[, , "CLUSTER"], at = c(0, 0.5, 1.5), col.regions = c("white",
+ "gray")))
```

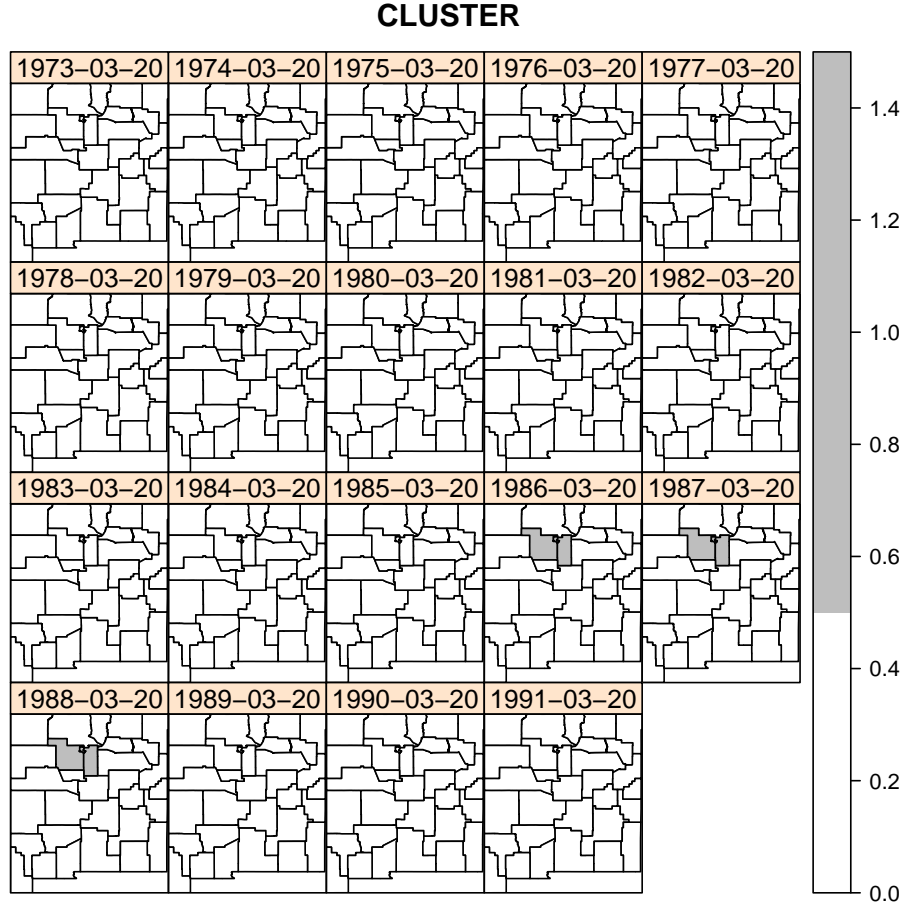


Fig:NMcluster}

Figure 5: Spatio-temporal cluster of brain cancer detected in New Mexico.

## 4. Zero-inflated models for cluster detection

[sec:zeroinfl}

Gómez-Rubio and López-Quílez (2010) extend this method to account for zero-inflation. In this case the observed number of cases come from a mixture distribution:

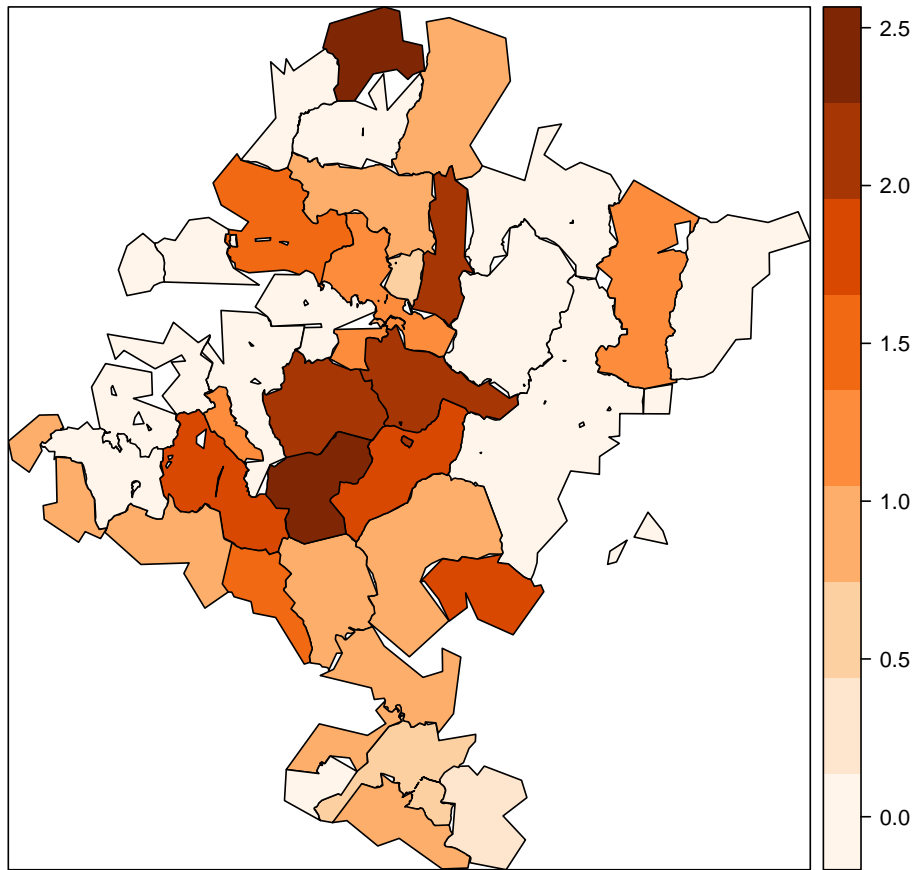
$$Pr(O_i = n_i) = \begin{cases} \pi_i + (1 - \pi_i)Po(0|\theta_i E_i) & n_i = 0 \\ (1 - \pi_i)Po(n_i|\theta_i E_i) & n_i = 1, 2, \dots \end{cases}$$

The relative risk  $\theta_i$  can be modelled using a log-linear model to depend on some relevant risk factors. Also, it is common that all  $\pi_i$ 's are taken equal to a single value  $\pi$ .

### 4.1. Brain Cancer in Navarre (Spain)

Ugarte, Ibáñez, and Militino (2006) analyse the incidence of brain cancer in Navarre (Spain). The aggregation level is the health district. Figure 4.1 shows the SMR. As it can be seen there are many areas where the SMR is zero because there are no cases in those areas. Ugarte,

Ibáñez, and Militino (2004) also tested for positive zero-inflation of these data compared to a Poisson distribution. The method implemented in this package is similar to the one used in Gómez-Rubio and López-Quílez (2010) for the detection of disease clusters of rare diseases.



{fig:Navarre}

Figure 6: SMR of brain cancer in Navarre (Spain).

## 4.2. Cluster detection

### *Cluster detection with no covariates*

Before starting our cluster detection methods, we will check the appropriateness of a Poisson GLM for this data. Fitting a log-linear model (with no covariates) gives the following model:

```
> m0 <- glm(OBSERVED ~ offset(log(EXPECTED)) + 1, family = "poisson",
+ data = brainnav)
> summary(m0)
```

Call:

```
glm(formula = OBSERVED ~ offset(log(EXPECTED)) + 1, family = "poisson",
```

```

data = brainnav)

Deviance Residuals:
    Min       1Q   Median       3Q      Max
-2.5227  -1.4783  -0.3203   0.7042   1.6393

Coefficients:
              Estimate Std. Error z value Pr(>|z|)
(Intercept) -7.752e-06  8.805e-02      0      1

(Dispersion parameter for poisson family taken to be 1)

    Null deviance: 63.733  on 39  degrees of freedom
Residual deviance: 63.733  on 39  degrees of freedom
AIC: 145.02

Number of Fisher Scoring iterations: 5

Furthermore, a quasipoisson model has been fit in order to asses any extra-variation in the
data:

> m0q <- glm(OBSERVED ~ offset(log(EXPECTED)) + 1, family = "quasipoisson",
+           data = brainnav)
> summary(m0q)

Call:
glm(formula = OBSERVED ~ offset(log(EXPECTED)) + 1, family = "quasipoisson",
    data = brainnav)

Deviance Residuals:
    Min       1Q   Median       3Q      Max
-2.5227  -1.4783  -0.3203   0.7042   1.6393

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept) -7.752e-06  9.703e-02      0      1

(Dispersion parameter for quasipoisson family taken to be 1.214555)

    Null deviance: 63.733  on 39  degrees of freedom
Residual deviance: 63.733  on 39  degrees of freedom
AIC: NA

Number of Fisher Scoring iterations: 5

The dispersion parameter in the previous model seems to be higher than 1, which may mean
that the Poisson distribution is not appropriate.

```

For this reason, and following Ugarte *et al.* (2004), a zero-inflated Poisson model has been fit. Here is the resulting model:

```
> m0zip <- zeroinfl(OBSERVED ~ offset(log(EXPECTED)) + 1 | 1, data = brainnav,
+   dist = "poisson", x = TRUE)
> summary(m0zip)
```

Call:

```
zeroinfl(formula = OBSERVED ~ offset(log(EXPECTED)) + 1 | 1, data = brainnav,
  dist = "poisson", x = TRUE)
```

Pearson residuals:

| Min     | 1Q      | Median  | 3Q     | Max    |
|---------|---------|---------|--------|--------|
| -1.3585 | -0.9137 | -0.1378 | 0.7137 | 1.8091 |

Count model coefficients (poisson with log link):

|             | Estimate | Std. Error | z value | Pr(> z ) |
|-------------|----------|------------|---------|----------|
| (Intercept) | 0.09347  | 0.09459    | 0.988   | 0.323    |

Zero-inflation model coefficients (binomial with logit link):

|             | Estimate | Std. Error | z value | Pr(> z ) |
|-------------|----------|------------|---------|----------|
| (Intercept) | -1.6158  | 0.6435     | -2.511  | 0.012 *  |

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Number of iterations in BFGS optimization: 9

Log-likelihood: -69.08 on 2 Df

Hence, the zero-inflated Poisson model will be used now to detect clusters of disease:

```
> brainnav$Expected <- brainnav$EXPECTED
> brainnavst <- STFDF(as(brainnav, "SpatialPolygons"), xts(1, as.Date("1990-01-01")),
+   brainnav@data, endTime = as.POSIXct(strptime(c("1990-01-01"),
+   "%Y-%m-%d"), tz = "GMT"))
> cl0 <- DetectClustersModel(brainnavst, coordinates(brainnav),
+   fractpop = 0.25, alpha = 0.05, typeCluster = "S", R = NULL,
+   numCPUS = 4, model0 = m0zip)
```

Library spdep loaded.

Library splancs loaded.

Library spacetime loaded.

Library DCluster loaded.

Library pscl loaded.

Library DClusterm loaded.

```
[1] 1 1
```

```
> cl0
```



|    | x        | y          | size | minDateCluster      | maxDateCluster      | statistic |
|----|----------|------------|------|---------------------|---------------------|-----------|
| 31 | 596886.8 | 4710520    | 4    | 1990-01-01 01:00:00 | 1990-01-01 01:00:00 | 2.520092  |
| 30 | 611795.5 | 4713762    | 3    | 1990-01-01 01:00:00 | 1990-01-01 01:00:00 | 2.016942  |
|    | cluster  | pvalue     |      |                     |                     |           |
| 31 | TRUE     | 0.02476587 |      |                     |                     |           |
| 30 | TRUE     | 0.04459518 |      |                     |                     |           |

As it can be seen, two clusters (with a p-value lower than 0.05) are detected. However, they overlap and we will just consider the one with the lowest p-value, which is shown in Figure 4.2.1

```
> names(cl0)[3] <- "size"
> knbinary(brainnav, cl0)
```

|    | CL1 | CL2 |
|----|-----|-----|
| 1  | 0   | 0   |
| 2  | 0   | 0   |
| 3  | 0   | 0   |
| 4  | 0   | 0   |
| 5  | 0   | 0   |
| 6  | 0   | 0   |
| 7  | 0   | 0   |
| 8  | 0   | 0   |
| 9  | 0   | 0   |
| 10 | 0   | 0   |
| 11 | 0   | 0   |
| 12 | 0   | 0   |
| 13 | 0   | 0   |
| 14 | 0   | 0   |
| 15 | 0   | 0   |
| 16 | 0   | 0   |
| 17 | 1   | 0   |
| 18 | 0   | 1   |
| 19 | 0   | 0   |
| 20 | 0   | 0   |
| 21 | 0   | 0   |
| 22 | 0   | 0   |
| 23 | 0   | 0   |
| 24 | 1   | 0   |
| 25 | 0   | 0   |
| 26 | 0   | 0   |
| 27 | 0   | 0   |
| 28 | 0   | 0   |
| 29 | 0   | 0   |
| 30 | 1   | 1   |
| 31 | 1   | 1   |
| 32 | 0   | 0   |

```

33  0  0
34  0  0
35  0  0
36  0  0
37  0  0
38  0  0
39  0  0
40  0  0

```

```

> brainnav$CLUSTER <- as.factor(knbinary(brainnav, c10)[, 1])
> levels(brainnav$CLUSTER) <- c("", "CLUSTER")

```

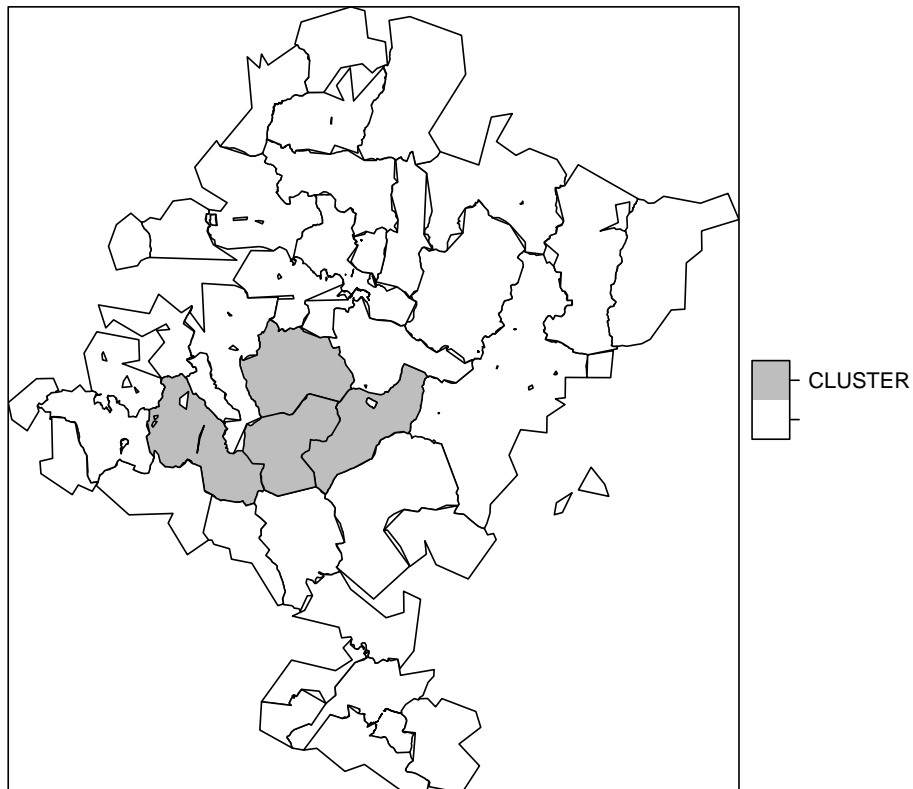


Fig:Navarrecl}

Figure 7: Cluster of brain cancer detected in Navarre (Spain).

{sec:mixed}

## 5. Mixed-effects models for cluster detection

{sec:bivar}

## 6. Bivariate models for cluster detection

{sec:disc}

## 7. Discussion

## References

- Bilancia M, Demarinis G (2014). “Bayesian scanning of spatial disease rates with integrated nested Laplace approximation (INLA).” *Statistical Methods & Applications*, **23**(1), 71–94. ISSN 1618-2510. doi:10.1007/s10260-013-0241-8. URL <http://dx.doi.org/10.1007/s10260-013-0241-8>.
- Gómez-Rubio V, López-Quílez A (2010). “Statistical methods for the geographical analysis of rare diseases.” *Advances in experimental medicine and biology*, **686**, 151–171.
- Gómez-Rubio V, Moraga P, Molitor J (2015). “Fast Bayesian classification for disease mapping and the detection of disease clusters.” *Submitted for publication*.
- Jung I (2009). “A generalized linear models approach to spatial scan statistics for covariate adjustment.” *Statistics in Medicine*, **28**(7), 1131–1143.
- Kulldorff M (1997). “A Spatial Scan Statistic.” *Communications in Statistics — Theory and Methods*, **26**(6), 1481–1496.
- Ugarte MD, Ibáñez B, Militino AF (2004). “Testing for Poisson Zero Inflation in Disease Mapping.” *Biometrical Journal*, **46**(5), 526–539.
- Ugarte MD, Ibáñez B, Militino AF (2006). “Modelling risks in disease mapping.” *Statistical Methods in Medical Research*, **15**, 21–35.
- Waller L, Turnbull B, Clark L, Nasca P (1992). “Chronic disease surveillance and testing of clustering of disease and exposure: application to leukemia incidence in TCE-contaminated dumpsites in upstate New York.” *Environmetrics*, **3**, 281–300.
- Waller LA, Gotway CA (2004). *Applied Spatial Statistics for Public Health Data*. John Wiley & Sons, Hoboken, New Jersey.
- Zhang T, Lin G (2009). “Spatial scan statistics in loglinear models.” *Computational Statistics and Data Analysis*, **53**(8), 2851–2858.

**Affiliation:**

Virgilio Gómez-Rubio  
Department of Mathematics  
School of Industrial Engineering  
University of Castilla-La Mancha  
02071 Albacete, Spain  
E-mail: [virgilio.gomez@uclm.es](mailto:virgilio.gomez@uclm.es)  
URL: <http://www.uclm.es/profesorado/vgomez>

Paula Moraga-Serrano  
London School of Hygiene & Tropical Medicine  
Keppel Street  
WC1E 7HT London, United Kingdom  
E-mail: [paula.moraga-serrano@lshtm.ac.uk](mailto:paula.moraga-serrano@lshtm.ac.uk)  
URL: <http://www.lshtm.ac.uk/aboutus/people/moraga-serrano.paula>

John Molitor  
College of Public Health and Human Sciences  
Oregon State University  
Corvallis, Oregon 97331, United States  
E-mail: [John.Molitor@oregonstate.edu](mailto:John.Molitor@oregonstate.edu)  
URL: <http://health.oregonstate.edu/people/molitor-john>

Barry Rowlingson  
Lancaster Medical School  
Furness Building  
Lancaster University  
Bailrigg, Lancaster LA1 4YG, United Kingdom  
E-mail: [b.rowlingson@lancaster.ac.uk](mailto:b.rowlingson@lancaster.ac.uk)  
URL: <http://www.lancaster.ac.uk/fhm/about-us/people/barry-rowlingson>