



## DClusterm: Model-based detection of disease clusters

Virgilio Gómez-Rubio

Universidad de  
Castilla-La Mancha

Paula Moraga-Serrano

London School of Hygiene & Tropical Medicine

---

### Abstract

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

---

## 1. Introduction

Kulldorff (1997) proposes a test for detecting disease clusters which will find the most likely cluster. This is called the Spatial Scan Statistic and the significance of the test is found via a Monte Carlo test. The test statistic is based on a likelihood ratio test for the following test:

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

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.

## 2. Generalised Linear Models for cluster detection

Jung (2009); Zhang and Lin (2009) 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 cluster whose coefficient is higher than 0 (i.e., increased risk).

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, so that

### 3. Spatio-temporal clusters

#### 3.1. Brain Cancer in New Mexico

The `brainNM` data set contains yearly cases of brain cancer in New Mexico from 1973 to 1991 (inclusive). 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.

```
> library(DClusterm)
> library(snowfall)
> 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. The SMR's have been plotted in Figure 3.1.

#### 3.2. Cluster detection

*Cluster detection with no covariates*

Similarly as in the spatial case, a GLM

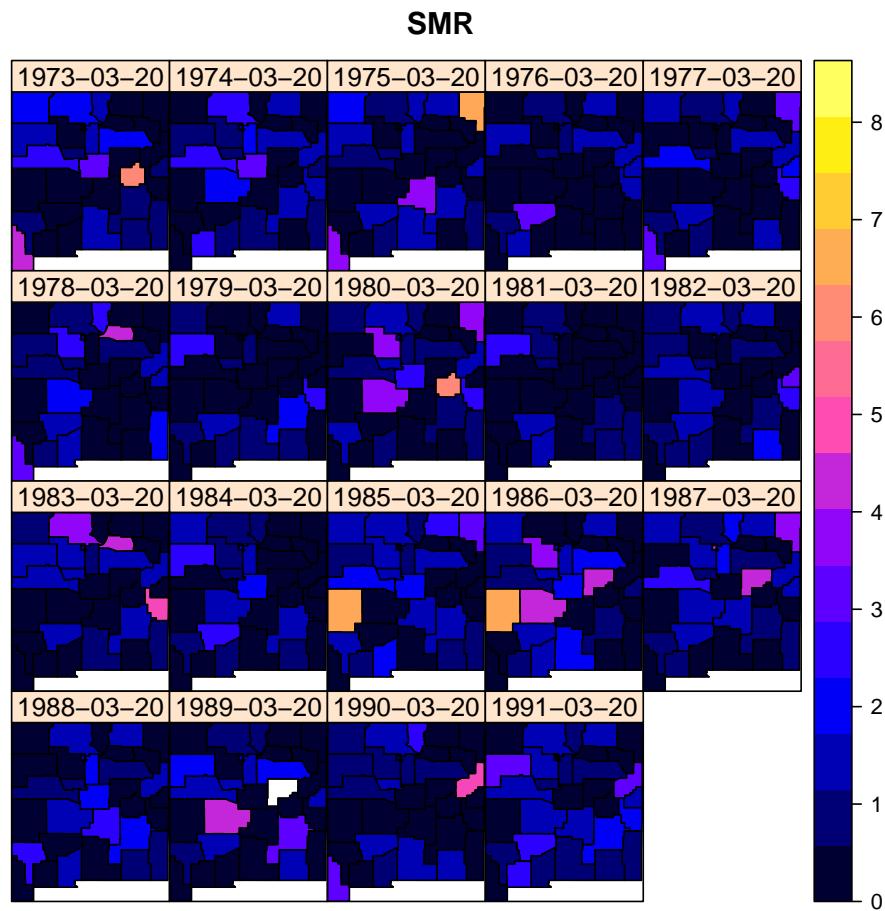


Figure 1: SMR of brain cancer in New Mexico.

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

Call:

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

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)	2.761e-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

```

```

> cl0 <- DetectClustersModel(brainst, coordinates(brainst@sp),
+   minDateUser = "1985-01-01", maxDateUser = "1989-01-01", fractpop = 0.15,
+   alpha = 0.05, typeCluster = "ST", R = NULL, numCPUS = 2,
+   model0 = m0)

```

```

> nrow(cl0)

```

```

[1] 180

```

```

> cl0[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 = 2, 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)))
```

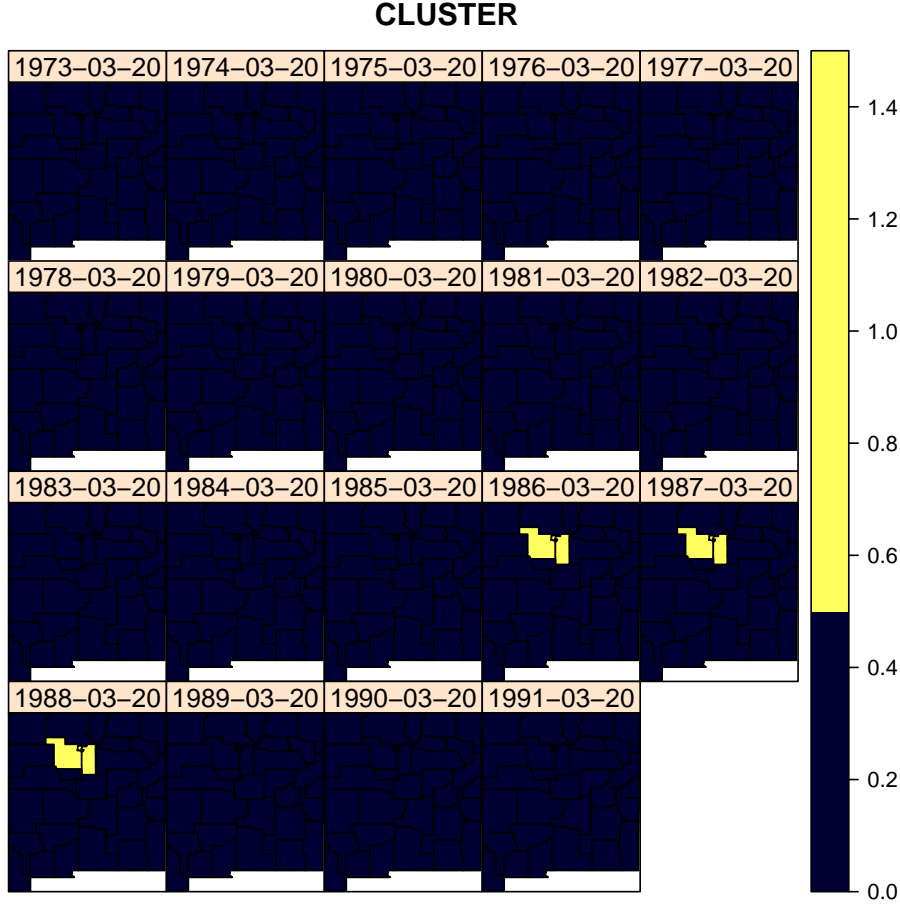


Figure 2: Spatio-temporal cluster of brain cancer detected in New Mexico.

#### 4. Zero-inflated models for cluster detection

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  $\pi$ .

#### References

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.

- 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.
- 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

Paula Moraga-Serrano  
London School of Hygiene & Tropical Medicine  
Keppel Street  
WC1E 7HT London, United Kingdom