1            **Supplementary material 1: statical analysis**

2

3    In this work, we use a statistical model for spatial data to make a prediction about what would have happened

4    in the year 2020 on the incidence of the disease, this is done since the results are greatly affected by the

5    pandemic. The statistical models for spatial data are divided by Cressie [1] into two broad classes:

6    geostatistical models with continuous spatial support and models in a lattice, also called area models [2], where

7    the data occur in a (possibly irregular) grid, with an enumerable set of vertices or locations . The two most

8    common area models are conditional autoregressive (CAR) models and simultaneous autoregressive (SAR)

9    models. These autoregressive models are used in many fields, including mapping disease rates [3], agriculture

10    [4], econometrics [5], ecology [6] and image analysis [7]. In this paper we will focus on CAR models. CAR

11    models are an example of the Gaussian Markov random fields [8] and the popular nested Laplace approach

12    integrated methods [9].

13

14    The basis of these models is the Gaussian Markov Fields. Random fields are multivariate distributions that are

15    generally used to describe the spatial association between variables X. A Markov random field extends the

16    Markov chain concept to a spatial context and assumes that such a joint distribution of X satisfies:

$$f(X_i | \boldsymbol{X}_{-i}) = f(X_i | \boldsymbol{X}_{j \sim i}),$$

17

18    where, $\boldsymbol{X}_{j \sim i}$ is the vector formed by all the components of X that are neighbors of i. A Gaussian Random

19    Markov Field (GMRF) is a Markov field where the random vector distribution (finite-dimensional) is a normal

20    or Gaussian distribution satisfying the conditional independence assumptions.

21

22    An n-dimensional random vector $y_{n \times 1} = (y_1, y_2, ..., y_n)^T$, $n < \infty$ has a n−variable distribution with mean vector $\mu_{n \times 1}$

23    and covariance matrix $\Sigma_{n \times n}$, and its probability density function (fdp) assumes the as follows:

24

$$f_{\boldsymbol{y}}(\boldsymbol{y}) = (2\pi)^{-n/2} |\boldsymbol{\Sigma}|^{-1/2} \exp\{-\frac{1}{2}(\boldsymbol{y} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{y} - \boldsymbol{\mu})\}.$$

25

26

27    This distribution will be denoted by $\boldsymbol{y} \sim N(\boldsymbol{\mu}, \Sigma)$ where $\mu$ and $\Sigma$ only have:

28

$$\mu_i = E(y_i)$$

30 and
$$\Sigma_{ij} = Cov(y_i, y_j), \ \Sigma_{ii} = Var(y_i) \ \text{and} \ Corr(y_i, y_j) = \Sigma_{ij}(\Sigma_{ii}\Sigma_{jj})^{-1/2}.$$

31

32 To build a GMRF we consider a graph G = (V, E) with n vertices where each vertex represents one of the

33 components of the vector y = $(y_1, y_2, ..., y_n)$ and edges connect nodes that have some sort of association. A GMRF

34 assumes that $\boldsymbol{y} = (y_1, y_2, \ldots, y_n)^T \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ and that the edges of the graph connect nodes i and j if

35 and only if $y_i \perp y_j | y_{-ij}$, that is, if $y_i$ is independent of $y_j$, given the components of y except $y_i$ and $y_j$. In a GMRF,

36 the covariance matrix brings information about the connections between the nodes through the precision matrix

37 $\Sigma^{-1} = Q$ which is a symmetric matrix and positive definite.

38

39 So, a random vector y = $(y_1, y_2, ..., y_n)^T \in R^n$ is called GMRF corresponding to a graph G = (V,E) with mean μ

40 and precision matrix Q > 0, if and only if the pdf of y has the following form:

41

42
$$\pi(\boldsymbol{y}) = (2\pi)^{-n/2}|\boldsymbol{Q}|^{1/2}\exp\left(-\frac{1}{2}(\boldsymbol{y} - \boldsymbol{\mu})^T\boldsymbol{Q}(\boldsymbol{y} - \boldsymbol{\mu})\right),$$

43

44 where the array Q satisfies the following condition:

45
$$Q_{ij} \neq 0 \Leftrightarrow \{i, j\} \in \mathcal{E}, \forall i \neq j.$$

46

47 If Q is a completely dense matrix, then G is fully connected, that is, the vertex is connected to all other vertices

48 in the graph. Let's focus on the case where Q is sparse. All results valid for the normal distribution will also

49 be valid for a GMRF. A detailed discussion of GMRF can be found in Rue and Held [8].

50

51 An example of a GMRF is the conditional autoregressive model or CAR model, in this case we consider a

52 geographic region that is partitioned into n subregions indexed by integers 1,2,...,n and assume that this

2

53    collection of sub-regions has a neighborhood system $\{V_i : i \ 1,...,n\}$, where $V_i$ denotes the collection of sub-

54    regions that, in a well-defined sense, are neighbors of the subregion i. In geographical terms,

55

56         $V_i = \{j :$ the subregions i and j share a boundary$\}$, to $i \in \{1,2,...,n\}$,

57

58    The neighborhood system is a key point in autoregressive or CAR models that are commonly used in spatial

59    statistics, the graphs that support the construction of the GMRF will be those that express these neighborhood

60    structures. In this context, the edges E in the graph G = (V, E), represent the connections in the geographic

61    structure and, consequently, define the neighbors that are used to model spatial dependence. The components

62    of the vector y are nodes of the graph.

63

64    Let $y_1,...,y_n$ be the observations made in the areas 1,...,n. Let us denote by $j \sim i$ that node j is a neighbor of node

65    i. The term conditional, in the CAR model is used because each element of the random process is conditionally

66    specified in the values of neighboring nodes, the CAR model assumes that the complete conditional

67    distributions are normal distributions. Then, we assuming that,

68

$$y_i | y_{-i} \sim N(\mu_i + \rho_{\mathcal{G}} \overline{(y - \mu)}_i, \frac{\sigma_{\mathcal{G}}^2}{d_i^{\mathcal{G}}})$$

69                                                                              ,

70

71    where $\sigma_{\mathcal{G}}^2/d_i^{\mathcal{G}}$ is the conditional variance of $y_i|y_{-i}$, $\rho_{\mathcal{G}}$ is a proportionality constant, $d_i^{\mathcal{G}}$ is the number of

72    neighbors of node i in the graph G, the average of the neighbors of node i is the:

$$\overline{(y - \mu)}_i = \sum_{\mathcal{E}^{\mathcal{G}}} (d_i^{\mathcal{G}})^{-1}(y_j - \mu_j)$$

73

74    and,

$$\mathcal{E}^{\mathcal{G}} = \{(i,j) \in E(\mathcal{G}) : j \sim i\}$$

75

76    is the set of edges that belong to the graph G. Consider the adjacency matrix $A_{\mathcal{G}} = (a_{ij})$ such that $a_{ii} = 0$, $a_{ij} = 1$

77    if i and j $a_{ij} = 0$ if i 6=$\sim$ j and $M_{\mathcal{G}} = diag\{d_1^{\mathcal{G}}, d_2^{\mathcal{G}}, \ldots, d_n^{\mathcal{G}}\}$.

Supplemental material

BMJ Publishing Group Limited (BMJ) disclaims all liability and responsibility arising from any reliance placed on this supplemental material which has been supplied by the author(s)

*BMJ Open Ophth*

78

79 Besag uses Brook's lemma [2, 10] and shows that when the matrix $(M_{\mathcal{G}} - \rho_{\mathcal{G}} A_{\mathcal{G}})^{-1}$ is positive definite and

80 symmetric the joint distribution for y is:

81 $$\boldsymbol{y} \sim N(\boldsymbol{\mu}, (\Sigma_{CAR}^{\mathcal{G}})^{-1}),$$

82

83 where $(\Sigma_{CAR}^{\mathcal{G}})^{-1} = \sigma_{\mathcal{G}}^2 (M_{\mathcal{G}} - \rho_{\mathcal{G}} A_{\mathcal{G}})^{-1}$. For the covariance matrix to be positive definite, it

84 is necessary that $\rho_{\mathcal{G}} < \frac{1}{\lambda_1}$ where $\lambda_1$ is the smallest eigenvalue of the matrix $M_{\mathcal{G}}^{-1/2} A_{\mathcal{G}} M_{\mathcal{G}}^{-1/2}$ Banerjee et al

85 [2].

86

87 In conclusion, the CAR model approach visualizes the geographic domain as an undirected graph with a vertex

88 in each region and an edge between two vertices if the corresponding regions share a geographic edge. This

89 creates well-defined neighbors for each region, which are used to define the joint or conditional distribution.

90 The distribution will be the multivariate normal distribution. All analysis of the CAR model is concentrated

91 on the covariance matrix Σ, which is defined by the graph of the geographic domain and the parameter ρ.

92

93 **REFERENCES:**

94 1. Cressie NAC. Statistics for spatial data. Rev. ed. New York: Wiley; 1993.

95 2. Banerjee S, Carlin BP, Gelfand AE, Banerjee S. Hierarchical Modeling and Analysis for Spatial Data. 0

96 edition. Chapman and Hall/CRC; 2003.

97 3. Elliott P, Wartenberg D. Spatial Epidemiology: Current Approaches and Future Challenges. Environmental

98 Health Perspectives. 2004;112:998–1006.

99 4. Besag J, Higdon D. Bayesian analysis of agricultural field experiments. J Royal Statistical Soc B.

100 1999;61:691–746.

101 5. LeSage J, Pace RK. Introduction to Spatial Econometrics. 0 edition. Chapman and Hall/CRC; 2009.

4

102    6. Arslan O, Akyürek Ö. Spatial Modelling of Air Pollution from PM10 and SO2 concentrations during Winter

103    Season in Marmara Region (2013-2014). International Journal of Environment and Geoinformatics. 2018;5:1–

104    16.

105    7. Besag J. On the Statistical Analysis of Dirty Pictures. Journal of the Royal Statistical Society Series B

106    (Methodological). 1986;48:259–302.

107    8. Rue H, Held L. Gaussian Markov Random Fields. 0 edition. Chapman and Hall/CRC; 2005.

108    9. Rue H, Martino S, Chopin N. Approximate Bayesian inference for latent Gaussian models by using

109    integrated nested Laplace approximations. Journal of the Royal Statistical Society: Series B (Statistical

110    Methodology). 2009;71:319–92.

111    10. Besag J. Spatial Interaction and the Statistical Analysis of Lattice Systems. Journal of the Royal Statistical

112    Society Series B (Methodological). 1974;36:192–236.

113