



University of Glasgow

SOME EXAM DATE
SOME EXAM TIME
EXAMINATION FOR THE DEGREES OF M.Sci. AND M.Sc.
(SCIENCE)

STATISTICS

Spatial Statistics 5M

“Hand calculators with simple basic functions (log, exp, square root, etc.) may be used in examinations. No calculator which can store or display text or graphics may be used, and any student found using such will be reported to the Clerk of Senate”.

NOTE: Candidates should attempt all questions.

1. (i) Consider the geostatistical process $\{Z(\mathbf{s})|\mathbf{s} \in \mathcal{D}\}$.
 - (a) Define what it means for $Z(\mathbf{s})$ to be weakly stationary. Then separately define what it means for $Z(\mathbf{s})$ to be isotropic. **[3 MARKS]**
 - (b) Give an example of a data set that is likely to be anisotropic. Justify your answer. **[2 MARKS]**
- (ii) Consider a geostatistical process with random variables $Z(\mathbf{s}_1)$ and $Z(\mathbf{s}_2)$, where the spatial locations are $\mathbf{s}_1 = (1, 0)$ and $\mathbf{s}_2 = (2, 0)$. Now consider a prediction location $\mathbf{s}_0 = (3, 0)$.
 - (a) Letting d denote the distance between 2 points, define the exponential covariance function at distance d with range parameter ϕ , nugget τ^2 and partial sill σ^2 . **[2 MARKS]**

CONTINUED OVERLEAF/

- (b) Now suppose that $\phi = 1$, $\tau^2 = 0$ and $\sigma^2 = 1$. Calculate the covariance matrix $\Sigma(\phi, \tau^2, \sigma^2)_{2 \times 2}$ for the random variables $(Z(\mathbf{s}_1), Z(\mathbf{s}_2))$ resulting from the exponential covariance function defined above. [2 MARKS]

- (c) Assuming that $Z(\mathbf{s})$ is a stationary zero-mean process with realisations $z(\mathbf{s}_1) = -0.5$ and $Z(\mathbf{s}_2) = 0.5$, compute the ordinary Kriging predictor for $Z(\mathbf{s}_0)$. [6 MARKS]

Hint - It may help you to remember that

$$\begin{pmatrix} a & b \\ c & d \end{pmatrix}^{-1} = \frac{1}{ad - bc} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix}.$$

- (d) Calculate the variance of your ordinary Kriging prediction at $Z(\mathbf{s}_0)$, and hence compute a 95% prediction interval for $Z(\mathbf{s}_0)$. [3 MARKS]
- (e) Given the location of the prediction point \mathbf{s}_0 in relation to the two data points $(\mathbf{s}_1, \mathbf{s}_2)$, is it sensible to make a prediction at \mathbf{s}_0 using data $(z(\mathbf{s}_1), z(\mathbf{s}_2))$? Justify your answer. [2 MARKS]

2. (i) Consider an areal unit process $\mathbf{Z} = (Z(\mathbf{s}_1), \dots, Z(\mathbf{s}_n))$, with corresponding neighbourhood matrix \mathbf{W} .

- (a) Define the k -nearest neighbours method for specifying a neighbourhood matrix \mathbf{W} , and give two drawbacks of using this approach. [3 MARKS]
- (b) Give an example of a situation where the k -nearest neighbours method would be preferable to the commonly used 'sharing a common border' method. Justify your answer. [2 MARKS]
- (c) Define the Local Indicator of Spatial Association (LISA) based on Moran's I statistic for region i . What can it tell you about the data in region i ? [3 MARKS]
- (d) The LISA was computed for three areas and gave values of: $I_1 = -0.5$, $I_2 = 0.03$, $I_3 = 0.45$. Describe what these three values tell you about the spatial dependence at areas $(1, 2, 3)$. [3 MARKS]

- (ii) Consider the following areal unit process and model: $\mathbf{Z} = (Z(\mathbf{s}_1), \dots, Z(\mathbf{s}_n)) \sim N(0, \mathbf{Q}^{-1})$, where \mathbf{Q} is the precision matrix and hence $\mathbf{Q}^{-1} = \Sigma$ is the covariance matrix.

- (a) What does it mean for the dependence between $(Z(\mathbf{s}_i), Z(\mathbf{s}_j))$ if: (i) the ij th element of \mathbf{Q} equals zero; (ii) the ij th element of Σ equals zero? [3 MARKS]
- (b) Suppose now that \mathbf{Z} comes from the *proper conditional autoregressive model*. Define \mathbf{Q} for this model, making sure that you define each of the elements in the formula. [3 MARKS]

CONTINUED OVERLEAF/

- (c) Is the *proper conditional autoregressive model* weakly stationary? Justify your answer. [3 MARKS]

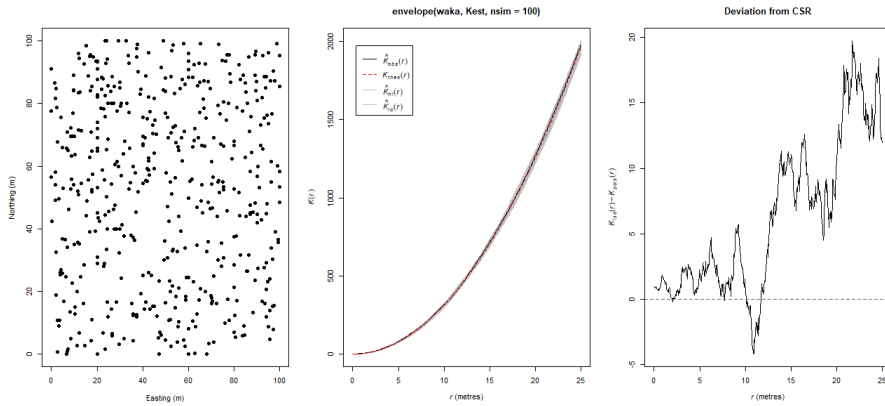
3. (i) Consider a spatial point process $Z = \{Z(A) : A \subset D\}$, where D is the domain of interest.

- (a) One hypothesis test for quantifying whether an observed spatial point pattern is completely spatially random is based on quadrat counts. Write down the null and alternative hypotheses for this test, the test statistic, and the distribution of the test statistic under the null hypothesis. [4 MARKS]

- (b) Consider an observed spatial point pattern with $n = 100$ points across a rectangular domain D . The rectangular domain is then split into 6 quadrats defined by two rows and three columns. The number of points in each of the six quadrats are: 20, 15, 10, 30, 12, 13. Conduct a test of complete spatial randomness and evaluate what it tells you about whether the observed point pattern is completely spatially random. [4 MARKS]

- (c) Give two downsides of the hypothesis test based on quadrat counts that you have just conducted. [2 MARKS]

- (ii) The figure below displays the locations of trees in a 100 metre square of the Waka national park in Gabon (left), as well as an estimate of Ripley's K function (middle) , and its deviation from complete spatial randomness (right).



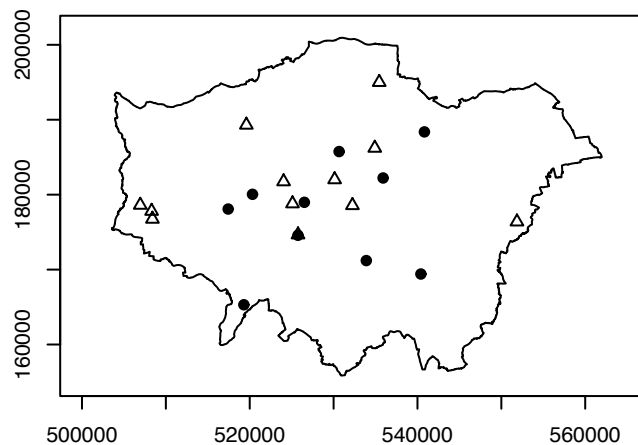
- (a) Define Ripley's K function and describe how it can be used to assess if an observed point pattern is completely spatially random. [2 MARKS]
- (b) From the figure above, do the trees appear to be completely spatially random? Justify your answer. [3 MARKS]

CONTINUED OVERLEAF/

(c) You decide to fit an inhomogeneous Poisson process model to the Waka tree data to estimate the extent to which the first order intensity function $\lambda(\mathbf{s})$ varies spatially. Would a parametric or non-parametric model be best for estimating $\lambda(\mathbf{s})$ here? Justify your answer. If you think a parametric model is best, write down a sensible model given the figure above. [3 MARKS]

(d) Give two reasons why an ecologist would be interested in estimating $\lambda(\mathbf{s})$ for in Waka forest. [2 MARKS]

4. (i) A researcher is interested in predicting the spatial pattern of nitrogen dioxide (NO_2) air pollution across Greater London (area within the border in the figure below), and collects data from the set of 21 air pollution monitors shown in the Figure below. The filled dots represent pollution monitors located at main roads, while the triangles represent monitors located in parks and other green spaces. It is known that NO_2 levels are higher near main roads and lower in green spaces.



(a) The aim of the analysis is to predict the spatial pattern of NO_2 levels across Greater London, and to do that the researcher choose to predict levels on a regular grid of prediction locations $\{\mathbf{s}_1^*, \dots, \mathbf{s}_N^*\}$. They propose the following analysis plan for analysing their data.

1. Plot the NO_2 data spatially to assess the presence of long-range trends.
2. Compute the empirical semi variogram for the raw NO_2 data with Monte-Carlo envelopes to assess the presence of spatial autocorrelation.
3. Fit a weakly stationary and isotropic geostatistical model to the data, with an exponential covariance function.

CONTINUED OVERLEAF/

4. Compare the fit of the model to the data with the fit of alternative plausible modelling choices for the covariance structure, such as replacing the exponential covariance function with a spherical covariance function.
5. Use the best fitting model (as chosen by AIC) to predict the NO_2 levels on the regular prediction grid $\{\mathbf{s}_1^*, \dots, \mathbf{s}_N^*\}$ using ordinary Kriging.

Provide a bullet pointed critique of this analysis plan, describing 3 weaknesses and 2 strengths of the proposed analysis. For each weakness you identify briefly describe how the analysis plan should be altered to address it. **[9 MARKS]**

- (b) The researcher asks your advice as to what results from their analysis they should present in their paper they plan to submit for publication. Write down 5 different outputs they should present in their paper to summarise their results. **[5 MARKS]**

(ii) Recall that the weakly stationary and isotropic Matérn covariance function is given by

$$\mathcal{C}(h) = \begin{cases} \sigma^2 + \tau^2, & h = 0, \\ \sigma^2 \frac{(h/\phi)^\xi}{2^{\xi-1}\Gamma(\xi)} K_\xi(h/\phi), & h > 0, \end{cases}$$

for distance h , where $\Gamma(\cdot)$ is the gamma function and $K_\xi(\cdot)$ is a modified Bessel function of the second kind. This function has 4 covariance parameters, the partial sill σ^2 , the nugget effect τ^2 , a range parameter ϕ and a smoothness parameter ξ .

- (a) What problem might occur when fitting this model and estimating $(\sigma^2, \tau^2, \phi, \xi)$ for a real data set? **[3 MARKS]**
- (b) How can this problem be overcome if one wishes to use the Matérn covariance function for their data? Ensure you provide a specific suggestion to overcome the problem identified above. **[3 MARKS]**

Total: 80

END OF QUESTION PAPER.