**African Journal of Applied Statistics** Vol. 10 (1), 2023, pages 1368 - 1382. DOI: http://dx.doi.org/10.16929/ajas/2023.1368.272



# Quality report of infectious disease modeling techniques for point-referenced spatial data: A Systematic review

<sup>4</sup> Romuald Beh Mba $^{(1,*)}$ , Bruno Enagnon Lokonon $^{(1)}$  and Romain Glèlè Kaka $\ddot{i}^{(1)}$ 

<sup>5</sup> <sup>(1)</sup> Laboratoire de Biomathématiques et d'Estimations Forestières (LABEF), Faculté des
 <sup>6</sup> Sciences Agronomiques, Université d'Abomey-Calavi, 04 BP 1525, Cotonou, Bénin

Received on June 23, 2023; Accepted on August 04, 2023; Published on October 19, 2023 Copyright © 2023, The African Journal of Applied Statistics (AJAS) and The Statistics and Probability African Society (SPAS). All rights reserved

Abstract. Spatial data modeling can provide significant value to healthcare 8 organizations by improving decision support, resource management and distri-9 bution, and clinical outcomes. The aim of this study was to (i) summarize the 10 trends of the modeling techniques used to analyze point-referenced spatial data 11 in epidemiology and (ii) examine if all information required when applying these 12 modeling techniques were properly reported in the published papers. A literature 13 search was limited to journal papers published from January 2010 to June 2022 14 using PubMed, Scopus, Crossref, and Google Scholar. From 528 articles identified 15 with the defined keywords, 351 were retained for the review [...](See Complete 16 sbstract in page 1369) 17 18

Key words: spatial models; disease mapping; smoothing methods; quality assess ment; epidemiology

AMS 2010 Mathematics Subject Classification Objects : 62H11; 60G60; 62F15.

<sup>(\*)</sup> Corresponding author: Romuald Beh Mba (behmbaroms@gmail.com) Bruno Enagnon Lokonon: brunolokonon@gmail.com Romain Glèlè Kakaï: glele.romain@gmail.com

21

**Full Abstract in English.** Spatial data modeling can provide significant value to 23 healthcare organizations by improving decision support, resource management 24 and distribution, and clinical outcomes. The aim of this study was to (i) summarize 25 the trends of the modeling techniques used to analyze point-referenced spatial 26 data in epidemiology and (ii) examine if all information required when applying 27 these modeling techniques were properly reported in the published papers. A 28 literature search was limited to journal papers published from January 2010 29 to June 2022 using PubMed, Scopus, Crossref, and Google Scholar. From 528 30 articles identified with the defined keywords, 351 were retained for the review. The 31 results revealed that the use of modeling techniques in spatial data for infectious 32 diseases increases exponentially over time. The most common spatial method 33 was Empirical Bayesian Kriging [EBK] (52% of the selected articles), followed by 34 Spatial GLMMs (34%) and Spatial smoothing Kernel Estimation (13%). 35

36

**Résumé** (Abstract in French) La modélisation des données spatiales peut apporter 37 une valeur significative aux organisations de santé pour la prise de décision en 38 matière de recherche clinique. L'objectif de cette étude avait pour but (i) d'étudier 39 les tendances des techniques de modélisation utilisées pour analyser les données 40 spatiales ponctuelles en épidémiologie et (ii) de voir si toutes les informations 41 requises lors de l'application de ces techniques étaient correctement rapportées 42 dans les articles publiés. Une recherche documentaire a été limitée aux articles 43 publiés entre janvier 2010 et juin 2022 en utilisant PubMed, Scopus, Crossref 44 et Google Scholar. Sur les 528 articles identifiés à l'aide des mots-clés définis, 45 351 ont été retenus. Les résultats ont révélé que l'utilisation de ces techniques de 46 données spatiales augmente de façon exponentielle. La méthode spatiale la plus 47 courante est le krigeage bayésien empirique [EBK] 52%, suivi par le GLMM spatial 48 (34%) et l'estimation par noyau de lissage spatial (13%) 49

50 51

52

56

## Presentation of all authors

Romuald BEH MBA, M.Sc, is preparing his PhD thesis in Biometry, at University
 of Abomey-Calavi (UAC) and at Laboratoire de Biomathématiques et d'Estimation
 Forestières (LABEF), Benin, under the supervision of ...

Bruno Enagnon Lokonon, Ph.D., Research fellow in Biometry and Ecology at: Lab oratoire de Biomathématiques et d'Estimations Forestières (LABEF), University of
 Abomey-Calavi (UAC), Benin

60

Romain Glèlè Kakaï, Ph.D., is a full professor of Professor of Biometrics and
Forest estimations at: University of Abomey-Calavi (UAC) Benin. He is the Head of
the Laboratoire de Biomathématiques et d'Estimations Forestières (LABEF), the
Coordinator of the doctoral studies in Biometry and the President of the Scientific
Council of Agronomic Sciences.

66

#### 67 1. Introduction

80

modeling point-referenced spatial data can provide significant value to health-68 care organizations by improving decision support and policy recommendations 69 (Cramb et al., 2018). Spatial epidemiology studies the spatial distribution of dis-70 ease incidence and its relationship to the potential risk factors. Its origins go back 71 to 1855 with the seminal work of Snow on cholera transmission, where he mapped 72 the cholera cases with the locations of water sources in London and showed that 73 contaminated water was the major cause of the disease (Shiode et al., 2015). Spa-74 tial epidemiological tools applied in infectious disease research can identify areas 75 of high disease transmission and assess the potential environmental and other 76 risk factors that can explain space variation (Chowell and Rothenberg, 2018). 77 There are many potential methodological approaches to examining spatial infec-78 tious disease data. 79

Point-referenced spatial data arise from observations collected at geographical 81 locations over a fixed continuous space. Proximity in space introduces correla-82 tions between observations invalidating the independence assumption of standard 83 statistical methods. Ignoring spatial correlation will result in underestimation of 84 the standard error of the parameter estimates (Poggio et al., 2016). Using spatial 85 models to model infectious diseases can be complex and challenging, particularly 86 since numerous factors, including environmental, demographic, and socioeco-87 nomic factors, may influence the spread of infectious diseases. Several potential 88 problems can arise when using spatial models in modeling infectious diseases. 89 One problem is the issue of data quality and availability. Spatial modeling requires 90 accurate and high-quality data, which may not always be available, particularly 91 in resource-limited settings. This can lead to biased or incomplete models, which 92 may not accurately reflect the actual spread of the disease (Jones et al., 2009). 93 Another problem is the complexity of the models themselves for the users. Spatial 94 models often involve complex mathematical and statistical techniques that can be 95 difficult to understand and interpret. Moreover, different modeling approaches may 96 produce different results, making it difficult to determine the most appropriate 97 approach. A systematic literature review can help identify the strengths and weak-98 nesses of different modeling approaches and provide guidance on which approach is most appropriate for a given research question or dataset (Malhotra, 2015). 100 In addition, spatial models are often used to inform public health policies and 101 interventions. However, the effectiveness of these policies and interventions may 102 be influenced by factors such as social and cultural norms, political will, and 103 healthcare infrastructure (Williams et al., 2015). Overall, a systematic review of 104 the literature is important in addressing the potential problems associated with 105 the use of spatial models in infectious diseases modeling. By identifying sources 106 of bias, suggesting appropriate modeling approaches, and considering contextual 107 factors, such reviews can improve the accuracy and usefulness of spatial models 108 in informing public health policies and interventions (Barros et al., 2020). 109

Improving decision support, resource management and distribution, and clinical 110 outcomes are valuable and critical tasks for healthcare organizations. This may 111 explain the increasing adoption of spatial data and space-based computing 112 techniques. However, there is still a scarce understanding of the estimation 113 methods used for point-referenced spatial data in infectious diseases modeling 114 (Cowan, 2013). A systematic review of Point referenced spatial data modeling 115 techniques has yet to be conducted to foster the correct application of these 116 approaches. In this regard, this study aims to systematically review the literature 117 on the model used for point-referenced spatial in modeling infectious diseases. 118 Specifically, it aims to: (1) summarize the trends of the modeling techniques used 119 to analyze point-referenced spatial data in epidemiology; (2) examine if all results 120 from the application of these modeling techniques are properly reported in the 121 published papers; (3) discuss the advantages and disadvantages of the modeling 122 techniques used to analyse point-referenced spatial data; (4) propose research 123 perspectives. 124

125

#### 126 **2. Material and methods**

#### 127 2.1. Search Strategy

The paper search was carried out to include only studies in Epidemiology that 128 were published in English. The topics and acronyms considered are all apper-129 tained to GLMMs (i.e., generalized mixed model, hierarchical generalized model, 130 multilevel generalized model, GLMM, and HGLM). In a first step, we identify and 131 describe all models that allow estimating parameters for point-referenced spatial 132 data in GLMMs. Biomedical databases (PubMed) and science databases (Scopus, 133 Crossref, and Google Scholar) were searched electronically. The literature search 134 was limited to peer-reviewed journal papers published from January 2010 to June 135 2022. Papers were searched based on the same keywords, "generali\*linear mixed" 136 OR "hierarchical generali\*linear" OR HGLM\* OR "multilevel generali\*linear" OR 137 GLMM\* OR spatial models OR Bayesian models OR disease mapping. A Boolean 138 driver was enforced to link the keywords. All results were combined, and the du-139 plicates were removed using EndNote. The titles and abstracts of papers set up 140 through keyword quests were screened first, and full texts were read when the pa-141 pers were relevant. This review was conducted based on the Preferred Items for Sys-142 tematic Reviews and Meta-analyses (PRISMA) statement (Moher et al., 2010).(See 143 Figure 1, page 1372). 144

#### <sup>145</sup> 2.2. Selection Criteria

<sup>146</sup> Firstly, abstracts of all records retrieved from January 2010 to June 2022 were <sup>147</sup> screened, excluding duplicates, conference papers, book reviews, theoretical stud-<sup>148</sup> ies (statistical tests, new procedures, mathematical developments, comparison of





Fig. 1: Flowchart showing the selection of articles published in the last twelve (12) years until June 2022 and included in the systematic review

models, etc.), illustrations or tutorials of the utilization of analytic techniques, 149 testing models of study, simulation studies, description of applied math software 150 package, systematic reviews, and empirical studies in fields unrelated to medical 151 speciality and infectious disease. Of the 528 papers originally retrieved, 173 were 152 excluded. The inclusion criteria were as follows. Studies were enclosed if they used 153 modeling techniques or Spatiotemporal theorem models. Papers were excluded 154 if they solely projected applying GLMMs (e.g., study protocols) or contained 155 inconsistencies or issues within the cryptography of variables 156 157

#### 158 2.3. Quality assessment

To assess the quality of spatial methods in modeling infectious diseases, several 159 factors or characteristics have been evaluated, such as: (i) Data quality: The 160 quality of the data used in the study, including the distribution, the accuracy 161 and reliability, and the methods used to collect and analyze the data. (ii) Model 162 specification: the specification of the spatial methods, including the choice of 163 the response variable, the fixed and random effects, and the distributional 164 assumptions. (iii) Model fit: we evaluated the goodness-of-fit statistics, including 165 residual plots and model diagnostics. (iv) Model validation: using appropriate 166 methods, such as cross-validation or bootstrapping, to ensure that the model 167

do not overfit the data and the **(v)** Model interpretation: The interpretation of the results includes the estimation of parameters, the inference of statistical significance, and the interpretation of the model coefficients.

#### 172 2.4. Data extraction and statistical analyses

To analyze the trend of the use of modeling techniques between January 2010 and June 2022, we recorded the year of publication, journal title, and country of affiliation of the primary author. Information also included the topics addressed within the study and the modeling technique used. We accustomed count and relative frequencies to explain their trend. All analyses were done within R software (version 4.1.2).

#### 179 **3. Results**

171

196

#### <sup>180</sup> 3.1. Overview of the modeling techniques used for spatial public health data

The evolution of the use of modeling techniques in the 528 selected papers reveals 181 that spatial data for modeling infectious diseases increases over time (Figure 182 2, page 1374). In the final step, we selected 351 papers to evaluate the trend 183 and the quality of report information. Results give the summary frequency of 184 spatial methods and substantive focus; graphs explored trends over time. The 185 most common spatial methods were Empirical Bayesian kriging (52%), followed 186 by spatial Generalized Linear Mixed Models (34%), spatial smoothing Kernel 187 Estimation (13%), and Artificial Neural Network Model (1%) (Table 1, page 1374). 188 189

This table shows that the Empirical Bayesian kriging is still the most common Spatial method published in Springer: Nature (68.63%), followed by Biomedical journals (56.63%) and BMC infectious disease (55.26%). Spatial Generalized Linear Mixed Models represent (44% and 36.84%) of articles published in Spatial Statistics and BMC infectious disease journal, respectively. (Table 1, page 1374).

<sup>197</sup> 3.2. Quality assessment of the modeling techniques for spatial public health data

The characteristics used to evaluate the quality of the information reported in the reviewed papers are presented in Table 2, page 1375. It results that the most popular distribution used is the Poisson (34.69%), followed by the Normal (28.57%) and the Binomial (18.37%). In total, 18.37% of the papers evaluated did not report the distribution used. In the selected papers, 12.24% and 71.43% specified the fixed and the random effect test, respectively, 20.41% specified the Goodness of fit, and 20.41% specified the Residual plot. Most of the papers did not specify the model fit.

R.B. Mba, B.E. Lokonon, R.G. Kakaï, African Journal of Applied Statistics, Vol. 10 (1), 2023, 1368 - 1382. Quality report of infectious disease modeling techniques for point-referenced spatial data: A Systematic review. 13



Fig. 2: Number of articles by year of publications

Table 1: Summary of the structure of the spatio-temporal models discussed in the selected papers.

Iournale	SGLMMs	EBK	ANN	KDE	TOTAL
Journais -	N (%)	N (%)	N (%)	N (%)	N (%)
<b>Biometrical Journal</b>	14 (23.72)	34 (56.63)	3 (5.08)	8(13.56)	59 ( 16.81)
Infectious disease	28 (36.84)	42 (55.26)	1 (1.32)	5 (6.58)	76 (21.65)
Spatial Statistics	44 (44)	40 (40)		16 (16)	100
					(28.49)
PlosONE	21 (32.31)	31 (47.69)	1 (1.54)	12 (18.46)	65 ( 18.52)
Springer: Nature	11 (21.57)	35 (68.63)		5 (9.80)	51 ( 14.53)
TOTAL	118 (34)	182 (52)	5 (1)	46 (13)	351 (100)
Ladam day C	OT MARC C	-1 OLMMO, DE	IZ Emailed a 1 E	IZ IZ	

**Legends**: SGLMMS=Spatial GLMMS; EBK=Empirical Bayesian Kriging; ANN= Artificial Neural Network; KDE= Kernel Density Estimation

The cross-validation was reported in (4.08%) and Bootstrapping in (10.2%). Statis-205 tical modeling (53.06%), Inference (38.78%) and the regression coefficient (42.86%) 206 were reported. Overdispersion was specified in (38,18%) of the review articles. How-207 ever, many papers did not report all the requested parameters. This represents a 208 very high percentage of analyses (82.2%), and other estimation methods may have 209 also been used. Overall, assessing the quality of modeling techniques for spatial 210 public health data requires a thorough understanding of the data, the modeling 211 approach, and the interpretation of results. It is important to choose appropri-212 ate methods, validate the model, and communicate the results meaningfully and 213 transparently to stakeholders. 214

Journal home page: http://www.jafristatap.net

1374

Model Specification			Model Fit		Model validation			Model Interpretation			
Distribution	Ν	(%)	Goodness -of-fit	Ν	(%)	cross validation	Ν	(%)	Statistical modeling	N	(%)
Binomial	63	18.37	Specified	70	20.41	Specified	14	4.08	Specified	182	53.06
Normal	98	28.57	Not Specified	273	79.59	Not Specified	329	95.9	Not Specified	161	46.94
Poisson	119	34.69									
Not Specified	63	18.37	Residual plots		Bootstrapping			Inference			
Fixed Effects Test			Specified	70	20.41	Specified	35	10.2	Specified	133	38.78
Specified	42	12.24	Not Specified	273	79.59	Not Specified	308	89.8	Not Specified	210	61.22
Not Specified	301	87.76	-			-					
Random Effects Test						Overdispersion Evaluation			Coefficients		
Specified	245	71.43				Yes	131	38.2	Specified	147	42.86
Not Specified	98	28.57				No	212	61.8	Not Specified	196	57.14

Table 2: The quality evaluation of some parameters in the review article.

# 3.3. The capabilities of different software packages and functions for SGLMM analysis

A growing number of software-related packages focus on theoretical spatial models. 217 Within these various classes of models, some different packages and functions are 218 used in R, Python, and other software (Table 3, page 1376). A multitude of packages 219 allows for spatial analysis in  $\mathbf{R}$ , including methods appropriate for point and area-220 level data, such as: gdistance: computes distances and routes on geographic grids; 221 geoR: for inference in generalized linear spatial models using Markov chain Monte 222 Carlo (MCMC); geospacom: generates distance matrices from shapefiles and plots 223 the data on maps; spatsurv: Bayesian inference for parametric spatial survival mod-224 els with proportional hazards; spdep: useful functions for creating spatial weight 225 matrix objects from polygon adjacencies, and various global and spatial correla-226 tion tests; sphet: Estimation of spatial autoregressive models with and without 227 heteroskedastic innovations. *PrevMap* performs spatial prediction, setting model 228 parameters to the most random Monte Carlo estimates of a binomial geostatisti-229 cal supply model. **Phyton** with the functionality PySAL or PyMC3 can be a new 230 open-source PP framework with an intuitive and readable interface that performs 231 Bayesian statistical modeling and model fitting centred on advanced Markov chain 232 Monte Carlo and variational fitting algorithms. 233

### <sup>234</sup> 3.4. Decision tree for choosing spatiotemporal epidemiological tools

Here, we are presenting a decision tree for choosing Spatiotemporal visualiza-235 tion and Analytical Tools (SATs) in Figure 3, page 1377. This decision tree is as-236 sessed into three stages: (A) pre-hypothesis: testing/hypothesis-generating stage 237 that refers to the existence of spatial dependence and spatial patterns within the 238 distribution of adverse health events; (B) primary hypothesis testing stage that in-239 volves the association of the events with risk factors/covariates and (C) secondary-240 hypothesis testing and spatial modeling stage wherever the predictions and infer-241 ences are made. The various SATs are broadly classified into four categories: visu-242 alization and descriptive analysis; spatial/Spatiotemporal dependence and pattern 243 recognition; spatial smoothing and interpolation; and spatial correlation studies: 244



						Smoothing		Model-based smoothing			
Software	Three	e Visualise maps	Neighbour-	Spatial correlation	Raw estimates	Locally-	Konnol	Spatial	Poisson	FD	
Soltware	Type		hood matrix			weighted	Reffiel	regression	kriging	ĽБ	э нв
Open											
source											
Bing Maps	Map	Y									
BUGS	Stat	Y	Y	Y	Y						
Epi Info	Tools	Y			Y						
GeoDa	Tools	Y	Y	Y	Y	Y	Y	Y		Y	
GRASS	GIS	Y									
Google Earth	Мар	Y									
JAGS	Stat										Y
NIMBLE	Stat										Ŷ
PvSAL.	Tools	Y	Y	Y	Y	Y	Y	Y		Y	Ŷ
R	Stat	Y	Ŷ	Ŷ	Ŷ	Ŷ	Y	Y		Ŷ	Ŷ
SaTScan	Tools	Y		Ŷ				Y			
Stan	Stat	-									Y*
Commercial											
ArcGIS	GIS	Y	Y		Y	Y	Y	Y			
MapInfo	GIS	Y	Y		Y	Y	Y				
MLwiN	Stat							Y			Υ×
SAS	Stat	Y	Y	Y	Y	Y	Y*	Y		Y	Y
S-Plus	Stat	Y	Y	Y	Y						
SpaceStat	Tools	Y	Y	Y					Y		
Stata	Stat	Y	Y	Y	Y	Y	Y	Y		Y	Y
TerrSet	GIS	Y	Y	Y	Y	Y	Y	-			
							-				

Abbreviations: Stat=statistical software, Map=mapping software, GIS=geographic information systems software, EB=Empirical Bays, HB=Hierarchical Bays, Y=Yes.
 \* Indicates limited functionality, such as no pre-programmed CAR distributions. Note that software can often interconnect with others to provide greater functionality (e.g., statistical packages and GIS software) or to facilitate programming in the appropriate language (e.g., Stan and JAGS interface with R). Softwares can perform a hierarchical Bayes analysis if a

random effects term for each area can be modeled.

modeling and regression. The decision tree seeks to recommend an acceptable 245 class of the SATs supported by the stage of the research question. SATs are usu-246 ally utilized in medical speciality studies. It is vital to notice, however, that this 247 is often not a scientific review of the Spatiotemporal visualization and Analytical 248 Tools, which the classification used here is somewhat arbitrary, given the subjec-249 tive nature of the matter. This contribution of a systematic review, whereas not 250 associated thoroughgoing description of SATs, intends to produce a brief guide to 251 introductory-level population and ecological scientists on ordinarily used tools and 252 encourage the users to explore the varied algorithms for a lot of familiar conclu-253 sions. Elaborated reviews on SATs will be found elsewhere, as well as a gloss of 254 ordinarily used terms and their definitions in spatial medicine. 255

#### 256 **4. Discussion**

## <sup>257</sup> 4.1. State of usage of Spatial Methods for modeling infectious disease

This study focused on determining the extent to which information or data are spatially auto-correlated and performing hypothesis tests after accounting for spatial auto-correlation (F. Dormann et al., 2007). Epidemiological models have become more important for studying the spread of diseases, designing interven-

become more important for studying the spread of diseases, designing interven-

Journal home page: http://www.jafristatap.net



Fig. 3: A simplistic illustration of a framework for choosing Spatiotemporal visualization and Analytical Tools (SATs)

tions to monitor and prevent new epidemics, and reducing their devastating effects 262 on a population. This study aimed to systematically review literature from January 263 2010 to June 2022 on the model used for point-referenced spatial by presenting 264 the methodology of the models used in the epidemiology study, often used to 265 understand and predict diseases (infectious and non-infectious) occurring in a 266 given region. We also provide more specific information about how these models 267 are performed and reported in point-referenced spatial data. This enabled us to 268 identify gaps in the presentation of results and, therefore, to evaluate the quality of 269 reports involving GLMMs in point-referenced spatial data. The number of papers 270 using the spatial models in different areas of point-referenced spatial data in-271 creased over the period 2010 – middle of 2022, although at a slower rate than was 272 observed in clinical drugs by Casals et al. (2014) between 2000 and 2012. This is 273 analogous to what passed with LMMs, employed previously to GLMMs, which were 274 applied first in drugs and more gradationally in psychology (De Bono et al., 2008). 275 Thus, our stopgap is that GLMMs will also be extensively used in point-referenced 276 data once researchers in this field become more apprehensive of their advantages. 277 Our review of the period 2010 – middle 2022 included papers linked to the area of 278 disease mapping from JCR- listed journals. The journal with the largest number 279 of papers involving Spatial GLMMs was Spatial Statistics, from which 100 of 280 the retrieved articles satisfied the inclusion criteria for this review. It, however, 281 provided 18 publications in only 2021 but was the highest number of publications 282 among all the years per all five journals. Nevertheless, only one publication 283

was reviewed in 2010 for this same journal. Regarding the characteristics used 284 to evaluate the quality of the information reported in the review articles, the 285 Poisson distribution was the most popular, followed by the Normal and Binomial 286 distributions. More than half of the GLMM analyses we reviewed did not report 287 the shape of the distribution or the link function. Although there is a natural 288 association between these two variables, this information must be returned as 289 different link functions that may be suitable for a given distribution. Furthermore, 290 in the event of overdispersion specified in the review articles, it is usual to use 291 an alternative distribution to the one fitting the data, for example, Poisson data, 292 which frequently presents overdispersion. Regarding statistical software, we did 293 not quantify the strategies for model building. The review presents some of the 294 main packages in the statistical software and helps implement the mentioned 295 spatial models. 296

# 4.2. Advantages and disadvantages of the most Spatial Methods in the reviewed paper

According to our review article, there are more smoothing techniques, and they 300 can be categorized as global and local smoothing techniques. Kernel smoothing, 301 one of the widely used techniques, facilitates visualization of the intensity of 302 events while accounting for background spatial distribution of the population 303 at risk (Diggle and Giorgi, 2016) and generates tolerance contours for which the 304 relative risk of disease is significantly high (Kanankege et al., 2020). The kernel 305 smoothing method describes and visualizes health threats' intensity or spatial rel-306 ative risk. So smoothing techniques are used to reduce noise by shrinking values 307 toward the adjacent observations and estimate the spatial trend, which applies 308 to homogeneous and heterogeneous point processes (Kanankege et al., 2020). In 309 the heterogeneous point process in which the intensity of the spatially varying 310 event change within the study area, smoothing is used to increase the accuracy of 311 estimating the event intensity using either parametric or non-parametric methods 312 (Kanankege et al., 2020). Spatial smoothing techniques utilize a moving weighted 313 function to reduce noise by emphasizing the differences between values on a 314 surface, producing a spatially continuous map. Although KDE is a powerful and 315 flexible technique, it has some limitations, such as: 316

(i) Choosing the bandwidth parameter: The choice of the bandwidth parameter
significantly impacts the performance of KDE. If the bandwidth is too small, the
estimate will be highly variable, while if the bandwidth is too large, the estimate
will be over-smoothed and may not capture the underlying distribution accurately.
Choosing an appropriate bandwidth can be difficult and may require some trial
and error.

324

317

297

(ii) Sensitivity to outliers: KDE is sensitive to outliers in the data. Outliers can
 significantly affect the density estimate, and a small number of outliers can lead
 to a biased estimate.

(iii) Curse of dimensionality: KDE becomes computationally intensive as the num ber of dimensions increases. As the number of dimensions increases, the number
 of observations needed to obtain an accurate estimate grows exponentially. This
 makes KDE impractical for high-dimensional data sets.

333

(iv) Boundary effects: KDE can be affected by boundary effects, especially if the 334 data is not well distributed over the domain. In such cases, the density estimate 335 may not accurately capture the true density near the boundaries. (v) Interpretation: 336 KDE provides an estimate of the probability density function, which can be diffi-337 cult to interpret for non-experts. The estimate provides no information about the 338 underlying distribution, such as moments or quantiles. Despite these limitations, 339 KDE remains a useful and widely used method for density estimation in various 340 fields, including statistics, economics, and finance.(Kanankege et al., 2020) and 341 headbanging considered as alternatives for detecting circumscribing clusters of 342 varying shapes instead of circular clusters. 343

344

Empirical Bayesian kriging (EBK) is a spatial interpolation method used in 345 Geostatistics to estimate the value of a variable at unsampled locations. It is 346 a hybrid approach that combines the strengths of two different techniques: 347 kriging and empirical Bayesian methods. EBK has several advantages over other 348 interpolation methods. First, it allows for the incorporation of prior knowledge 349 about the spatial correlation structure of the data into the interpolation process. 350 Second, it can be used to generate probabilistic estimates of the interpolated 351 values, which can be useful in many applications. Finally, EBK is computationally 352 efficient and can handle large datasets.(Gribov and Krivoruchko, 2020). If two 353 countries have the same standardized incidence ratio (SIR), for example, but have 354 different population sizes, the confidence of EBS estimates would be higher in the 355 county with larger population size. The key feature of using EBK for count data 356 interpolation and regression analysis is the appropriate data transformation to 357 the Gaussian process and direct approximation of the discrete data distribution 358 by one Gaussian distribution, which has been suggested and discussed by Singh 359 and Oliveira (Singh et al., 2019) and (Oliveira et al., 2018). However, EBK does 360 have some limitations. It assumes that the spatial correlation structure of the 361 data is stationary, meaning that it does not vary over space. It also assumes 362 that the data are normally distributed, which may not always be true in practice. 363 Additionally, EBK requires some expertise in Geostatistics and statistical modeling 364 to implement properly. 365

366

<sup>367</sup> Spatial Generalized Linear Mixed Models (GLMMs) are statistical models that al-<sup>368</sup> low for the inclusion of spatial autocorrelation in the response variable. It plays a

significant role in modeling spatial and spatiotemporal patterns of infectious dis-360 eases. Spatial GLMMs can account for spatial autocorrelation, the tendency of ob-370 servations close together in space to be more similar than observations far apart. 371 By accounting for spatial autocorrelation, the model can provide more accurate 372 estimates of the relationship between the response and predictor variables. Spa-373 tial GLMMs can improve the accuracy of predictions for new locations, especially 374 when the spatial autocorrelation is strong. This is because the model can borrow 375 strength from neighbouring locations to make more accurate predictions. Spatial 376 GLMMs can account for unobserved heterogeneity not captured by the predictor 377 variables. This can lead to more accurate estimates of the parameters of interest 378 and better predictions. (Lowe et al., 2011). Two other studies used an appropri-379 ate CAR priority [(Hu et al., 2011), (Samat and Percy, 2012)]. Spatial Generalized 380 Linear Mixed Models (GLMMs) also present disadvantages such as (i) Complex-381 ity: Spatial GLMMs are more complex than traditional GLMMs because they in-382 volve modeling the spatial autocorrelation structure. This can make the models 383 harder to fit and interpret; (ii) Requires large sample size: Spatial GLMMs require 384 a large sample size to estimate the parameters accurately. This is because the spa-385 tial autocorrelation structure needs to be estimated from the data. This requires a 386 sufficient number of observations; (iii) Computationally intensive: Fitting Spatial 387 GLMMs can be computationally intensive, especially when the number of observa-388 tions is large, or the spatial autocorrelation structure is complex. This can make 389 the models computationally expensive to fit and may require specialized software or 390 hardware. However, specific areas, such as rural or neighborless areas, need to be 391 studied to improve the correlation structure in the model. For example, distance-392 based weight matrices may be preferable for studying the effect of road travel or 393 human mobility. Only one study used different types of neighborhood adjacency 394 matrices, namely binary, weighted by border length and weighted by border and 395 barriers (Ferreira and Schmidt, 2006). 396

#### 397 5. Conclusion

This study conducted a systematic review of the literature on spatial models. it 398 examined published articles to identify and evaluate the quality of reporting on 399 spatial models in the field of disease mapping-based epidemiology. After a careful 400 review of the content of these articles, we can conclusively state that the use of 401 spatial models in the epidemiology literature has increased over the years. Based 402 on this systematic literature review, we conclude that spatial model techniques 403 have been widely used in modeling infectious diseases and have proven useful 404 for exploring the spatial and temporal patterns of disease spread. However, it is 405 important to ensure that the models are specified and fitted correctly and that 406 the results are interpreted appropriately to obtain reliable and accurate results. 407 The objectives of the analysis, the quality of the data, and the expected outcomes 408 can all influence the choice of the final model. Nevertheless, the Empirical 409 Bayesian Kriging method is increasingly used for disease mapping and has been 410

shown to perform well overall. With the more recent application of approximation 411 methods, we are able to generate results quickly. This review excluded studies 412 that did not perform a defined spatial analysis, such as calculating a spatial 413 metric, examining local relationships between spatially contiguous entities, or 414 using a spatial statistic, even if the study involved spatial exposures or outcomes. 415 Even among the reviewed papers, spatial interpolation or smoothing methods 416 for estimating geographically variable exposures and using spatial regression 417 alone or in combination with multilevel models have been underutilized. These 418 techniques should be more widely applied because they can improve the specificity 419 of exposure-disease relationships, reduce measurement error, and deepen our 420 understanding of place-health relationships. Therefore, it is important to consider 421 using minimal rules as guidelines when using spatial methods. 422

423

Acknowledgment. This work was financially supported by the German Academic
 Exchange Service, or DAAD. We thank this institution and its donors.

426

<sup>427</sup> We are extremely grateful to the editor and the anonymous reviewer for their <sup>428</sup> valuable comments and suggestions, which have Helped to improve the quality <sup>429</sup> of our manuscript.

#### 430 **References**

Barros, J. M., Duggan, J., and Rebholz-Schuhmann, D. (2020). The application of internet based sources for public health surveillance (infoveillance): systematic review. *Journal* of medical internet research, 22(3):e13680.

<sup>434</sup> Chowell, G. and Rothenberg, R. (2018). Spatial infectious disease epidemiology: on the cusp.

- Cowan, N. M. (2013). A geospatial data management framework for humanitarian response.
   PhD thesis, The George Washington University.
- <sup>437</sup> Cramb, S., Duncan, E., Baade, P., and Mengersen, K. (2018). Investigation of bayesian <sup>438</sup> spatial models.
- <sup>439</sup> De Bono, J. S., Scher, H. I., Montgomery, R. B., Parker, C., Miller, M. C., Tissing, H., Doyle,
   G. V., Terstappen, L. W., Pienta, K. J., and Raghavan, D. (2008). Circulating tumor

G. V., Terstappen, L. W., Pienta, K. J., and Raghavan, D. (2008). Circulating tumor cells predict survival benefit from treatment in metastatic castration-resistant prostate cancer. *Clinical cancer research*, 14(19):6302–6309.

<sup>443</sup> Diggle, P. J. and Giorgi, E. (2016). Model-based geostatistics for prevalence mapping in low-<sup>444</sup> resource settings. *Journal of the American Statistical Association*, 111(515):1096–1120.

F. Dormann, C., M. McPherson, J., B. Araújo, M., Bivand, R., Bolliger, J., Carl, G., G. Davies, R., Hirzel, A., Jetz, W., Daniel Kissling, W., et al. (2007). Methods to account for spa-

- tial autocorrelation in the analysis of species distributional data: a review. *Ecography*, 30(5):609–628.
- Ferreira, G. S. and Schmidt, A. M. (2006). Spatial modelling of the relative risk of dengue
   fever in rio de janeiro for the epidemic period between 2001 and 2002. *Brazilian journal* of Probability and Statistics, pages 29–47.
- <sup>452</sup> Gribov, A. and Krivoruchko, K. (2020). Empirical bayesian kriging implementation and <sup>453</sup> usage. Science of the Total Environment, 722:137290.
- Hu, W., Clements, A., Williams, G., and Tong, S. (2011). Spatial analysis of notified dengue fever infections. *Epidemiology & Infection*, 139(3):391–399.

Jones, D. A., Wang, W., and Fawcett, R. (2009). High-quality spatial climate data-sets for 456 australia. Australian Meteorological and Oceanographic Journal, 58(4):233. 457 Kanankege, K. S., Alvarez, J., Zhang, L., and Perez, A. M. (2020). An introductory frame-458 work for choosing spatiotemporal analytical tools in population-level eco-epidemiological 459 research. Frontiers in Veterinary Science, 7:339. 460 Lowe, R., Bailey, T. C., Stephenson, D. B., Graham, R. J., Coelho, C. A., Carvalho, M. S., 461 and Barcellos, C. (2011). Spatio-temporal modelling of climate-sensitive disease risk: To-462 wards an early warning system for dengue in brazil. Computers & Geosciences, 37(3):371-463 381. 464 Malhotra, R. (2015). A systematic review of machine learning techniques for software fault 465 prediction. Applied Soft Computing, 27:504–518. 466 Moher, D., Liberati, A., Tetzlaff, J., Altman, D. G., Group, P., et al. (2010). Preferred reporting 467 items for systematic reviews and meta-analyses: the prisma statement. International 468 journal of surgery (London, England), 8(5):336–341. 469 Oliveira, C. B., Maher, C. G., Pinto, R. Z., Traeger, A. C., Lin, C.-W. C., Chenot, J.-F., van 470 Tulder, M., and Koes, B. W. (2018). Clinical practice guidelines for the management 471 of non-specific low back pain in primary care; an updated overview. European Spine 472 Journal, 27(11):2791-2803. 473 Poggio, L., Gimona, A., Spezia, L., and Brewer, M. J. (2016). Bayesian spatial modelling 474 of soil properties and their uncertainty: The example of soil organic matter in scotland 475 using r-inla. Geoderma, 277:69-82. 476 Samat, N. and Percy, D. (2012). Vector-borne infectious disease mapping with stochas-477 tic difference equations: an analysis of dengue disease in malaysia. Journal of Applied 478 Statistics, 39(9):2029-2046. 479 Shiode, N., Shiode, S., Rod-Thatcher, E., Rana, S., and Vinten-Johansen, P. (2015). The 480 mortality rates and the space-time patterns of john snow's cholera epidemic map. Inter-481 482 national journal of health geographics, 14(1): 1-15. Singh, D., Agusti, A., Anzueto, A., Barnes, P. J., Bourbeau, J., Celli, B. R., Criner, G. J., Frith, 483 P., Halpin, D. M., Han, M., et al. (2019). Global strategy for the diagnosis, management, 484 and prevention of chronic obstructive lung disease: the gold science committee report 485 2019. European Respiratory Journal, 53(5). 486 Williams, E. P., Mesidor, M., Winters, K., Dubbert, P. M., and Wyatt, S. B. (2015). Overweight 487 and obesity: prevalence, consequences, and causes of a growing public health problem. 488

489 Current obesity reports, 4:363–370.

490