# Practical considerations on data patterns in Bayesian Maximum Entropy Estimation: A systematic and critical review

**Emmanuel Ehnon Gongnet[a,b,*], Romaric Vihotogbé[c], Tranquillin Affossogbe Sédjro Antoine [a] and Romain Glèlè Kakaï[a]**

[a] Laboratoire de Biomathématiques et d'Estimations Forestières, Université d'Abomey-Calavi, Benin
[b] Institut Tchadien de Recherche Agronomique pour le Développement, Tchad
[c] École de Foresterie Tropicale, Université National d'Agriculture, Benin
*Corresponding author: Emmanuel EHNON GONGNET: ehnon.gongnet@gmail.com

## ABSTRACT

*Objective:* It is well known that some data features (sample size, skewness, among others) may determine method performance. The choice of those features depends on the researcher's level of awareness on the statistical method. In this study, the level of awareness on the influence of spatial data key characteristics (sample size, skewness, spatial dependency and variogram model) in Bayesian Maximum Entropy (BME) was analyzed.

*Methodology:* A systematic review was conducted that covers the period from 1990 (year of BME introduction) to 2019. Two main keywords "Bayesian Maximum Entropy" and "BME" were used for literature search. Publications which only mentioned the keywords without applying BME were excluded while those with application and/or BME theory discussion were considered. Six of the world's leading Open Access sources of scientific literature were considered, namely: Science Direct, African Journals Online, Springer, Google Scholar, MPDI and Academic Journals. A total of 118 research articles from 62 journals were identified. The sample sizes screened shows that 25.4% of the published articles used few samples (less than 100), which implies the variogram might not yield accurate results. The analysis of the use of skewness showed that most researchers do not apply transformation on skewed data (82.2%) nor consider skewness in their descriptive statistics (90.7%). Even though 11% of theoretical papers have mentioned about spatial dependency level, 92.4% of them failed to consider it. Most researchers (68.64%) do not specify the variogram models but when they do, they mostly use exponential model (12.7%). It clearly appears in this review that most researchers do not consider the effect of sample size, skewness, and spatial dependency level when applying BME. Yet very few research works have focused on these aspects. This therefore calls for more in-depth studies on the effect of data characteristics on BME's performance.

**Keywords:** Bayesian Maximum Entropy, sample size, skewness, spatial dependency.

## INTRODUCTION

In geostatistics spatial, temporal and or spatio-temporal data with various characteristics are used to predict values of a random variable at any unsampled geographical location. The traditional method for this prediction is the stochastic one called kriging (Mazari, 2012). Compared to other interpolation techniques, kriging has been demonstrated to be the best, since it quantifies, in addition to estimating associated uncertainty (Goovaerts, 2001; Adhikary & Dash, 2017). However, the kriging technique requires data to be at the intrinsic stationarity (Meul & Meirvenne, 2003). This guarantees the existence of the covariance and assumes that the observations variance exists and is only a function of the lag distance (Obaid & Mohammed, 2020). Moreover, similar with most statistical methods (e.g. regression analysis and analysis of variance), the accuracy of kriging estimates depends on the degree of skewness of the dataset (Arslan, 2017). Ignoring important assumptions, e.g. normality and homogeneity of variance can lead to either Type I or Type II error therefore affect the conclusion (Osborne, 2010). To overcome these weaknesses regarding the kriging techniques, Bayesian Maximum Entropy (BME) method was introduced (Christakos, 1990). The BME is a strong mathematical background – based approach for handling spatiotemporal data (Serre *et al.,* 1998; He & Kolovos, 2017). The strength of this technique relies on its ability to integrate various sources of data, regardless of their nature (Gengler & Bogaert, 2016). Therefore, BME method has proven to be more reliable for spatial and spatiotemporal analyses through a wide range of applications in environmental geology, environmental sciences, soil sciences, public health, ecology, remote sensing, energy, real estate research, among others (Yang *et al.,* 2016). However, its reliability may depend on data pattern and structure such as sample size, degree of skewness, spatial dependence and variogram model which might influence the prediction accuracy. In general, statistical analyses are less biased with high sample size (Steven *et al.,* 1998; Neerchal *et al.,* 2008; Serge & Brigittte, 2012). In geostatistics, the variogram summarizes variations of a variable in a targeted region (Lark, 2000). The computation of this central statistic in spatial data analysis is affected by the sample size (Christakos, 2000), large sample size yield better variograms (Webster, 1992) while smaller sample sizes ($<$ 50) lead to erratic variograms (Webster, 1992), which shows the strong effects of sample size on BME – based estimates. BME is more robust than other methods for spatiotemporal prediction. However, in terms of data normality handling (Christakos, 2000), the commonly used Matheron's variogram estimator (Matheron, 1963) in geostatistics is still based on variances and thus, is sensitive to data normality (Kerry & Oliver, 2007). Moreover, the measure of entropy is affected by the degree of skewness. Therefore, the higher the skewness is, the smaller the calculated entropy (Orton & Lark, 2009). These fluctuations may have significant effects on BME estimates, unfortunately, skewness was not addressed in many research (Fu *et al.,* 2014; Hosseini & Kerachian, 2017; Xiao *et al.,* 2018). In some cases, transformation is applied on skewed data (Lee & Ellis, 1997; Douik *et al.,* 2005; Jiang *et al.,* 2014), to improve the normality and therefore the accuracy of estimates (Amin *et al.,* 2018). Whether a statistical test is considered "robust" to violations of data normality or a nonparametric test, taking normality into account can improve the accuracy of the results (Osborne, 2010). Thus, if the variogram and entropy are sensitive to data skewness, it is important to consider how skewed data would be handled in various fields of application using BME method.

Spatial dependency captures information on autocorrelation among locations at a given lag

distance apart in a targeted geographical domain. It helps to unveil unobservable heterogeneity when data are sampled from large geographical areas. Spatial dependence is traditionally described using the variogram which is strongly influenced by the marginal distribution of the random field (Kazianka & Pilz, 2010). The Spatial Dependency level is estimated based on the nugget to sill ratio. There is a strong spatial dependency, if the ratio is less than 0.25 and moderate spatial dependency if the ratio ranges between 0.25 - 0.75 (El-Sayed Ewis, 2012). In soil sciences, strong spatially dependent characteristics may be influenced by intrinsic variations in soil properties such as texture and organic matter content. The moderately spatially dependent variables, for example bulk density and total porosity, are controlled more by extrinsic variations such as cultivation (Jerry and Sidney, 2012). In an ecological application, Jetz and Rahbek (2002) and Lichstein *et al*. (2002) demonstrated that spatial dependency influences the model's coefficient of determination. In order to design better geostatistical – based systems' control, this research built up a state of knowledge regarding how raw primary data characteristics are taken into account during geostatistical analysis, by addressing the following questions: (1) What are the sample sizes often used in BME analysis? (2) How is data skewness handled when applying BME analyses? (3) To what extent is spatial dependency considered?

## METHODOLOGY

The materials included original research articles, reviews, as well as letters to editors. Six of the worlds' leading Open Access sources of scientific publications were identified. These were Science Direct (www.sciencedirect.com), African Journals Online (www.ajol.info), Springer (www.springer.com), Google Scholar (www.scholar.google.com), MPDI (www.mpdi.com) and Academic journals (www.academicjournals.org). This review focused on research works published on BME from 1990 up to December 2019. The articles were downloaded, using the search keywords: "Bayesian Maximum Entropy" and/or "geostatistics". Only research with BME application and theorical articles with no application, were analysed. Then, basic information was recorded on journal name, impact factor, title, objectives, major findings, and keywords. Each paper was reviewed and information extracted on sample sizes, degree of skewness, option for data transformation, transformation method, spatial dependency level and strategy for handling it, and variogram model.
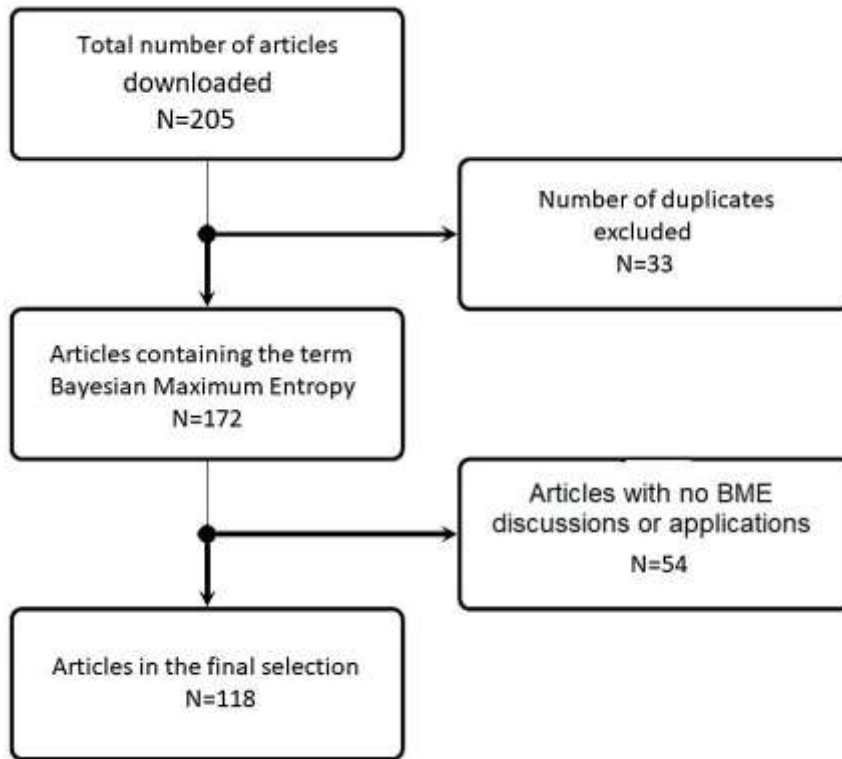
**Fig.1.** Article screening diagram flow

**Data analysis:** The frequencies of the considered sample size, skewed data, data transformation, spatial dependence and variogram models were computed. The diversity of BME application was obtained by computing the frequency per field. Based on this frequency, histograms were computed to easily visualize the importance per field of application. The BME evolution from 1990 to 2019 was described using a line plot combined with characteristic equation of the curve and coefficient of determination ($R^2$).

**Sample size:** Three classes of sample were used, that is small (less than 100), average (100-1000) and large (greater than 1000). Frequency of each class by fields of application were computed and the relationship between the sample classes and fields were tested using a Chi-square test ($x^2$) of independence.

**Skewness:** The distribution of datasets on which BME were applied in published articles were explored using descriptive statistics. The frequencies of the response to the question: "Is data transformed (Yes/No)?", "data transformation methods", and "Is there any descriptive statistics (Yes/No)?" were computed by field to build a contingency table. A Chi-square test ($x^2$) of independence was used to find out if the decision to transform data, computation of descriptive statistics and the choice of transformation methods was field dependent.

**Spatial dependence:** The frequency of papers in which spatial dependency is considered and the ones where it is not considered was computed. Chi-square test ($x^2$) of independence was used to establish the relationship between the consideration of spatial dependence and field of application.

**Variogram models:** Summary statistics of variograms models in the published articles were computed. Correlation between field and the choice of variogram model was assessed using Chi-square test ($x^2$) of independence.

## RESULTS

**Characteristics of published papers on BME:** One hundred and eighteen articles from 60 journals (Appendix 2) were downloaded and analysed. The BME were applied in many fields, with soil sciences (19.8 %), hydrogeology (15.3%) and health science (15.3%) as the major fields (Figure 2). The BME evolution plot from 1990 to 2019 shows that it is increasingly used with the average number following an exponential distribution and coefficient of determination ($R^2$) of 59%. This suggests that the model explained 59% of the variations. However, three time periods can easily be distinguished from the reviewed articles, viz: 1990 - 2000 with 3 publications per year, 2000-2010 with up to 7 per year and 2010-2019 with a sharp increase up to 15 article per year (Figure 3).
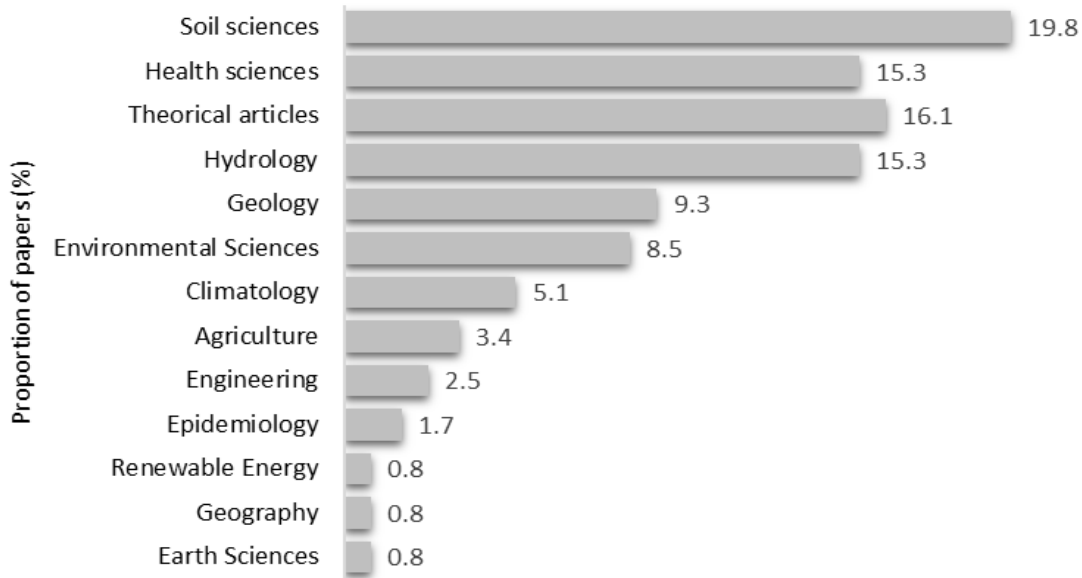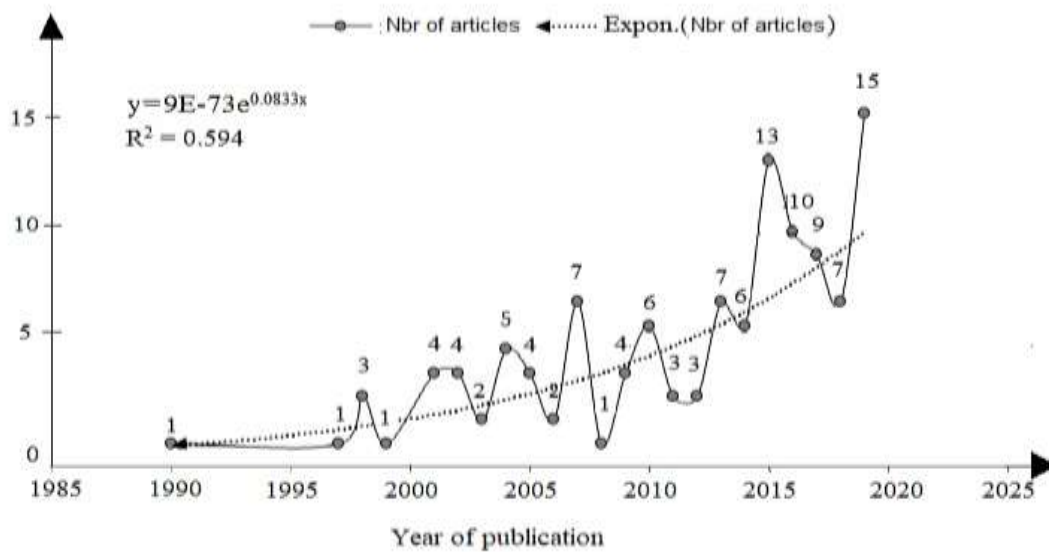


**Fig.2:** BME application fields



**Fig.3:** Evolution of BME–based research as calculated from publications

**Handling sample size in BME application:** Of all the sample sizes used by researchers, three classes were observed: small (less than 100), average (100-1000) and large (greater than 1000) sample sizes. Despite the diversity of sample sizes used, average sample size accounting to 32.2% was most frequent. But considering the field of application, the lower sample sizes were mainly used in engineering (66.7%), climatology (50%) and geology (45.5%). Average sample sizes were considered in earth sciences (100%), epidemiology (100%), geography (100%) and soil sciences (52.2%), while the larger samples were used in renewable energy (100%), agriculture (50%), climatology (50%), and environmental sciences (50%). Some authors (19.5%) did not specify the sample sizes in their articles (Table 1). In addition, the choice of sample size was not field dependent ($x^2$ (36,118) = 40.23, p=0.29) (Table 2).

**Table 1:** Proportion of publications (%) by sample size and field of application

| Fields | $\chi^2$ (36,118) =40.23, p=.29 | | | | |
|---|---|---|---|---|---|
| | **Unknown** | **<100** | **100 - 1000** | **> 1000** | |
| Agriculture | 25.0 | 25.0 | 0.0 | 50.0 | |
| Climatology | 0.0 | 50.0 | 0.0 | 50.0 | |
| Earth Sciences | 0.0 | 0.0 | 100.0 | 0.0 | |
| Engineering | 33.3 | 66.7 | 0.0 | 0.0 | |
| Environmental Sciences | 0.0 | 20.0 | 30.0 | 50.0 | |
| Epidemiology | 0.0 | 0.0 | 100.0 | 0.0 | |
| Geography | 0.0 | 0.0 | 100.0 | 0.0 | |
| Geology | 9.1 | 45.5 | 36.4 | 9.1 | |
| Health science | 33.3 | 27.8 | 27.8 | 11.1 | |
| Hydrology | 22.2 | 27.8 | 27.8 | 22.2 | |
| Remote sensing | 100.0 | 0.0 | 0.0 | 0.0 | |
| Renewable Energy | 0.0 | 0.0 | 0.0 | 100.0 | |
| Soil sciences | 4.3 | 13.0 | 52.2 | 30.4 | |
| Theoretical articles | 42.1 | 21.1 | 26.3 | 10.5 | |
| Total | 19.5 | 25.4 | 32.2 | 22.9 | |

**BME application on skewed data:** The BME was applied on data with all kind of characteristics, especially with regard to skewness. Skewness values on which BME was applied since 1999 ranges between −2.68 and 32.5, with a mean of 2.98. Hence, the majority of the BME - based analyses used positively skewed data. In addition, most attributes were highly peaked from a normal distribution (kurtosis = 3) with the kurtosis, ranging between −0.6 and 28.76 and mean value 5.44 (Table 2). Thus, when data is skewed, decision for transformation or to use normality parameters (skewness and kurtosis value) as descriptive statistics, did not depend on research field $\chi^2 (13, 118) = 14.33$, p = 0.35 and $\chi^2 (13,118) = 11.91$, p = 0.54, respectively (Table 3). However, cases of transformation accounted for 17.8%, out of which only 9.3% considered skewness value in their descriptive statistics (Table 3). Specifically, data transformation was applied in earth sciences (100%), geology (36.4%), soil sciences (26.1%), and agriculture (25%), while skewness was mostly used as descriptive statistics in soil sciences (26.1%), climatology (16.7%), and environmental sciences (10%) (Table 3). Two transformation technics were often applied in geostatistical analyses: the logarithmic and Box-Cox transformations (Table 4). However, none of the applications of these two techniques was field - dependent

$\chi^2(39,118) = 30.57$, p = 0.83). In general, logarithmic transformation was the most often used (6.8%), while Box-Cox represented only 0.8%, both representing 38% and 5% of transformations, respectively. Some authors (10.2%) had transformed data but failed to specify the method, which accounts for 57% transformations. The logarithmic transformation was mostly applied in earth sciences (100%), with moderate application in the other sciences. Box-Cox transformation was only applied in soil sciences (4.3%). In most cases (82.2%), researchers did not mention anything about data transformation methods considered.

**Table 2:** Degree of skewness and kurtosis in the reviewed articles

| Parameters | Skewness | Kurtosis |
|---|---|---|
| Minimum | -2.68 | -0.6 |
| Maximum | 32.5 | 28.76 |
| Standard deviation | 6.13 | 7.92 |
| Mean | 2.98 | 5.44 |

**Table 3**: Level of application of data transformation

| Fields | Is data transformed? | | Is there any descriptive statistics? | |
|---|---|---|---|---|
| | $\chi^2$ (13, 118) =14.33, p=.35 | | $\chi^2$ (13, 118) =11.91, p=.54 | |
| | **No (%)** | **Yes (%)** | **No (%)** | **Yes (%)** |
| Agriculture | 75.0 | 25.0 | 100.0 | 0.0 |
| Climatology | 100.0 | 0.0 | 83.3 | 16.7 |
| Earth Sciences | 0.0 | 100.0 | 100.0 | 0.0 |
| Engineering | 100.0 | 0.0 | 100.0 | 0.0 |
| Environmental Sciences | 80.0 | 20.0 | 90.0 | 10.0 |
| Epidemiology | 100.0 | 0.0 | 100.0 | 0.0 |
| Geography | 100.0 | 0.0 | 100.0 | 0.0 |
| Geology | 63.6 | 36.4 | 90.9 | 9.1 |
| Health science | 77.8 | 22.2 | 94.4 | 5.6 |
| Hydrology | 88.9 | 11.1 | 100.0 | 0.0 |
| Remote sensing | 100.0 | 0.0 | 100.0 | 0.0 |
| Renewable Energy | 100.0 | 0.0 | 100.0 | 0.0 |
| Soil sciences | 73.9 | 26.1 | 73.9 | 26.1 |
| Theoretical articles | 94.7 | 5.3 | 94.7 | 5.3 |
| **TOTAL** | **82.2** | **17.8** | **90.7** | **9.3** |

**Table 4:** Decision for data transformation before BME application in different fields

| Field | Transformation option $\chi^2$ (39,118) =30.57, p=.83 | | | |
|---|---|---|---|---|
| | Box–Cox | Logarithmic | Unknown | None |
| Agriculture | 0.0 | 0.0 | 25.0 | 75.0 |
| Climatology | 0.0 | 0.0 | 0.0 | 100.0 |
| Earth Sciences | 0.0 | 100.0 | 0.0 | 0.0 |
| Engineering | 0.0 | 0.0 | 0.0 | 100.0 |
| Environmental Sciences | 0.0 | 10.0 | 10.0 | 80.0 |
| Epidemiology | 0.0 | 0.0 | 0.0 | 100.0 |
| Geography | 0.0 | 0.0 | 0.0 | 100.0 |
| Geology | 0.0 | 9.1 | 27.3 | 63.6 |
| Health science | 0.0 | 11.1 | 11.1 | 77.8 |
| Hydrology | 0.0 | 0.0 | 11.1 | 88.9 |
| Remote sensing | 0.0 | 0.0 | 0.0 | 100.0 |
| Renewable Energy | 0.0 | 0.0 | 0.0 | 100.0 |
| Soil sciences | 4.3 | 8.7 | 13.0 | 73.9 |
| Theoretical articles | 0.0 | 5.3 | 0.0 | 94.7 |
| Total | 0.8 | 6.8 | 10.2 | 82.2 |

**Spatial dependency in BME application:** The results (Figure 4) shows that majority of researchers in all fields do not account for spatial dependency level (92.4%). It was observed that only the fields of Remote sensing (100%), Geology (36.4%) and Environmental Sciences (20%) considered the spatial dependency levels in their papers. In addition, results showed that the consideration of spatial dependency in a given paper depends on the field of BME application ($x^2$ (13, 118) = 33.76, p=.001).
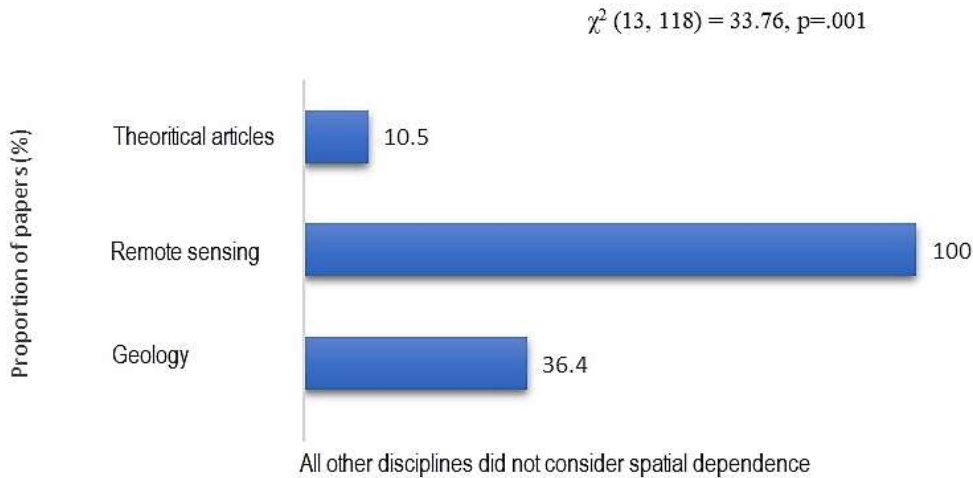
$\chi^2$ (13, 118) = 33.76, p=.001



**Figure 4:** Consideration of Spatial dependency

**Choice of variogram models in BME computations:** Two types of variogram models were involved in all geostatistical analyses: simple and nested variogram models. The most used models were exponential variogram (12.71%) and nugget, exponential,

and spherical models combined (4.24%) (Table 5). The choice of the variogram model is significantly linked to specific fields of research with $x^2$ (130,118) = 167.10, p = 0.02 (Table 6). In simple model, the exponential model was mostly used by climatologists while spherical model was used in engineering and geology. However, remote sensing, agriculture, hydrology and health science mostly used a combination of exponential, gaussian and spherical models.

**Table 5**: Variogram models in the published articles

| Type of Model | Variogram | Frequencies (%) |
|---|---|---|
| | Exponential | 12.71 |
| Simple | Gaussian | 2.54 |
| | Spheric | 3.39 |
| | Unspecified | 68.64 |
| | Nugget effect + exponential + exponential | 1.69 |
| | Nugget effect + gaussian + exponential | 1.69 |
| | Nugget effect + exponential | 1.69 |
| Nested | Nugget effect + spheric + gaussian | 0.85 |
| | Nugget effect + exponential + spherical | 4.24 |
| | Nugget effect + exponential + gaussian + spherical | 1.69 |
| | Nugget effect + spheric + holesin + nugget effect | 0.85 |
| | Total | 100 |

**Table 6**. Correlation between field and the choice of variogram model

| Fields | p-value = 0.02; $\chi^2$ = 167.10; df =130 (S ) | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | **1** | **2** | **3** | **4** | **5** | **6** | **7** | **8** | **9** | **10** | **11** |
| Agriculture | 0 | 0 | 25 | 0 | 0 | 0 | 75 | 0 | 0 | 0 | 0. 0 |
| Climatology | 0 | 0 | 0 | 0 | 16.7 | 0 | 66.7 | 0 | 0 | 0 | 16. 7 |
| Earth Sciences | 0 | 0 | 0 | 0 | 0 | 0 | 100 | 0 | 0 | 0 | 0. 0 |
| Engineering | 0 | 0 | 0 | 0 | 0 | 0 | 66.7 | 0 | 33.3 | 0 | 0. 0 |
| Environmental Sciences | 30 | 0 | 0 | 0 | 0 | 0 | 70 | 0 | 0 | 0 | 0. 0 |
| Epidemiology | 50 | 0 | 0 | 0 | 0 | 0 | 50 | 0 | 0 | 0 | 0. 0 |
| Geography | 0 | 0 | 0 | 0 | 0 | 0 | 100 | 0 | 0 | 0 | 0. 0 |
| Geology | 18.2 | 0 | 0 | 9.1 | 9.1 | 0 | 45.5 | 0 | 18.2 | 0 | 0. 0 |
| Health science | 0 | 0 | 0 | 5.6 | 0 | 0 | 88.9 | 0 | 0 | 5,6 | 0. 0 |
| Hydrology | 5.6 | 0 | 5.6 | 16.7 | 0 | 11.1 | 55.6 | 5.6 | 0 | 0 | 0. 0 |
| Remote sensing | 0 | 100 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0. 0 |
| Renewable Energy | 0 | 0 | 0 | 0 | 0 | 0 | 100 | 0 | 0 | 0 | 0. 0 |
| Soil sciences | 17.4 | 4.3 | 0 | 0 | 4.3 | 0 | 69.6 | 4.3 | 0 | 0 | 0. 0 |
| Theory | 21.1 | 0 | 0 | 0 | 0 | 0 | 73.7 | 0 | 5..3 | 0 | 0. 0 |

[1]Exponential, [2]Exponential + Exponential, [3]Exponential + Gaussian + Spherical, [4]Exponential + Spherical, [5]Gausian, [6]Gaussian + exponential, [7]NR, [8]Nugget + exponential, [9]Spheric, [10]Spheric + Gaussian, [11]Spheric + holesin + nugget.

## DISCUSSION

The review of the use of BME package within the time frame of 1990 to 2019 shows the increasingly wide application of this statistical tool in various fields (Jerry & Sidney, 2012; He & Kolovos, 2017). The exponential trends in the application of the BME is due to its superiority compared to other methods. In general, small sample sizes less than 100 were mostly used in BME application in every discipline, except in Earth Sciences, Geography, Epidemiology, and Renewable Energy. It is demonstrated that variogram estimation significantly depends on sample size, with 100 to 150 locations ensuring optimal variogram calculation (Webster and Oliver, 1992). In this review, 25.4% of sample sizes were small, indicating low reliability of the associated research, because of their low variogram performance (Lark *et al.,* 2017). However, Lark (2000) demonstrated that variogram performance also depends on estimators used. Therefore, an optimum of 60 and 90 – 120 locations are needed when applying maximum likelihood and method of moments, respectively, for variogram calculation (Lark, 2000). Our data showed that BME was mostly applied on positively skewed data. This suggests that attributes on which BME were applied were dominated by low values and the arithmetic may less describe the central tendency of datasets (Clay *et al.,* 1999). Indeed, natural variables are highly skewed. The kurtosis lies between -0.6 and 28.76 with a mean of 5.44 indicating that the attributes are highly peaked. This result shows that BME is applied on variables that have greater deviation from normal distribution. However, 82.2% of researchers did not mention whether they transformed data before applying BME despite the fact that error pdf's change considerably as the skewness values vary (Christakos, 2000). The variogram is sensitive to highly positive skewed data due to some exceptionally large values (Webster & Oliver, 2001). Kerry & Oliver (2007) have demonstrated using a simulated data that when the skewness value is outside the bounds of $\pm 1$, the variogram for the transformed data is more suitable than the variogram for the original data. He also found that when the skewness coefficient is large, the form of the experimental variogram becomes erratic and is difficult to model. Among authors that have transformed data (17.8%), logarithmic transformation was the most used (38%) and 57% of them failed to specify the transformation. However, Manikandan (2010) suggested that a method of transformation should be selected based on the relationship between the standard deviation and the mean. Logarithmic transformation should be used when the standard deviation approximates the mean (data is positively skewed). Square root transformation should be preferred when the mean is proportional to the variance. If the standard deviation approximates the mean squared, a reciprocal transformation can be performed. Box-Cox transformation, which represented only 5% transformations, covers all traditional methods (e.g., square root, log, inverse, cubic root) and easily produces optimal normalization (Osborne, 2010). In order to overcome the effect of highly heterogenous and skewed data, such as the distribution of monthly Haemorrhagic fever with renal syndrome (HFRS) cases, a class-dependent Bayesian Maximum Entropy (cd-BME) was introduced. The method demonstrated a greater capacity in modelling the variability in HFRS data by dividing the original dataset into discrete incidence classes (He *et al.,* 2019). The main objective of any geostatistical analysis is to predict attributes values at unsampled locations using the concept of random function, which assumes that, the set of unknown values are spatially dependent random variables. This existence of spatial dependence between observations is essential to mapping (Goovaerts, 1998). The spatial dependence can be measured by a

correlogram or semi variogram and classified into weak, moderate, or strong using the nugget-to sill ratio (Cambardella *et al.,* 1994). However, this study showed that most researchers do not consider the level of spatial dependency in their studies (92.7%) despite the fact that it can affect the variogram

computation and the BME prediction accuracy. The higher the spatial dependence, the higher the accuracy of the prediction (Orton & Lark, 2007; Orton & Lark, 2009). When the spatial dependency value increases, BME accuracy of estimation becomes less important (Zimmerman & Zimmerman, 1991).

## CONCLUSION

This review showed that most researchers involved in spatiotemporal prediction and mapping have been neglecting important factors such as skewness, sample size and spatial dependence which might influence BME accuracy. It is clear that when data is highly skewed, variogram becomes erratic and is difficult to model (Kerry & Oliver, 2007). In this review, 25.4% of articles published used small sample sizes that might not allow variogram to yield accurate results. Large samples are costly and time consuming, therefore recommending large sample to reach

optimal variogram estimation can be an obstacle to the large adoption of this method (Lark, 2000). Also, when sample sizes increase, BME computation becomes difficult especially for non-gaussian variables (Cao *et al.,* 2014). However, among the articles reviewed, there have been no study fully investigating empirically BME robustness. Therefore, there is need for an empirical evaluation of the effect of sample size, skewness, and spatial dependency structure on BME prediction.

## REFERENCES

Adhikary PP. and Dash C, 2017. Comparison of deterministic and stochastic methods to predict spatial variation of groundwater depth. Applied Water Science 7(1): 339-348. https://doi 10.1007/s13201-014-0249-8

https://doi.org/10.1007/s13201-014-0249-8

Amin A, Shah B, Khattak AM, Baker T, Anwar S, 2018. Just-in-time customer churn prediction: With and without data transformation. In 2018 IEEE congress on evolutionary computation (CEC) (pp. 1-6). IEEE.

https://doi.org/10.1109/CEC.2018.8477954 PMCid: PMC5755424

Arslan H, 2017. Determination of temporal and spatial variability of groundwater irrigation quality using geostatistical techniques on the coastal aquifer of Çarşamba Plain, Turkey, from 1990 to 2012. Environmental Earth Sciences 76(1): 1-12.

https://doi.org/10.1007/s12665-016-6375-x

Cambardella C, Moorman T, Parkin T, Karlen D, Turco R, Konopka A, 1994. Field scale variability of soil properties in Central Iowa soils. Soil Sci. Soc. Am. J. 58:1501- 1511.

https://doi.org/10.2136/sssaj1994.0361599500 5800050033x

Cao G, Yoo EH, Wang S, 2014. A statistical framework of data fusion for spatial prediction of categorical variables. Stochastic environmental research and risk assessment 28(7): 1785-1799.

https://doi.org/10.1007/s00477-013-0842-7

Christakos G, 1990. A Bayesian/maximum-entropy view to the spatial estimation problem. Mathematical Geology 22(7): 763-777.

https://doi.org/10.1007/BF00890661

Christakos G, 2000. Modern spatiotemporal geostatistics (Vol. 6). Oxford university press. Christakos G, 2012.

Random field models in earth sciences. Courier Corporation.

Clay A, Jason G, Forcella F, Ellsbury M, Carlson C,1999. Sampling weed spatial variability on the field wild scale. Weed Science 47:674-681.

https://doi.org/10.1017/S0043174500091323

El-Sayed Ewis O, 2012. Improving the prediction accuracy of soil mapping through geostatistics. International Journal of Geosciences vol. 2012.

Gengler S. and Bogaert P, 2016. Bayesian data fusion applied to soil drainage classes spatial mapping. Mathematical Geosciences 48(1): 79-88.

https://doi.org/10.1007/s11004-015-9585-y

Goovaerts P, 1998. Geostatistical tools for characterizing the spatial variability of microbiological and physico-chemical soil properties Biol Fertil Soils 27:315-334.

https://doi.org/10.1007/s003740050439

Goovaerts P, 2001. Geostatistical modelling of uncertainty in soil science. Geoderma 103: 3-26.

https://doi.org/10.1016/S0016-7061 (01)00067-2

He J, Christakos G, Wu J, Jankowski P, Langousis A, Wang Y, Zhang W, 2019. Probabilistic logic analysis of the highly heterogeneous spatiotemporal HFRS incidence distribution in Heilongjiang province (China) during 2005-2013. PLOS Neglected Tropical Diseases 13(1), e0007091. doi:10.1371/journal.pntd.0007091

https://doi.org/10.1371/journal.pntd.0007091

PMid: 30703095 PMCid: PMC6380603

He J. and Kolovos A, 2017. Bayesian maximum entropy approach and its applications: a review. Stoch. Environ. Res. Risk Assess 32(4): 859-877. https://doi.org/10.1007/s00477-017-1419-7

https://doi.org/10.1007/s00477-017-1419-7

Hosseini M. and Kerachian R, 2017. A Bayesian maximum entropy- based methodology for optimal spatio-temporal design of groundwater monitoring networks. Environmental Monitoring and Assessment 189:433

https://doi.org/10.1007/s10661-017-6129-6

PMid: 28779429

Hu J, Zhou J, Zhou G, Luo Y, Xu X, Li P, Liang J, 2016. Improving estimations of spatial distribution of soil respiration using the Bayesian maximum entropy algorithm and soil temperature as auxiliary data. Plos one, 11(1): e0146589.

https://doi.org/10.1371/journal.pone.0146589

PMid: 26807579 PMCid: PMC4726581

Jerry AN. and Sidney RV, 2012. Spatial patterns and correlation of soil properties of a lowland soil. Journal of Soil Science and Environmental Management 3(12):287-296.

Jetz W. and Rahbek C, 2002. Geographic range size and determinants of avian species richness. Science 297(5586):1548-1551.

https://doi.org/10.1126/science.1072779

PMid: 12202829

Kazianka H. and Pilz J, 2010. Copula-based geostatistical modelling of continuous and discrete data including covariates. Stochastic environmental research and risk assessment 24(5):661-673

https://doi.org/10.1007/s00477-009-0353-8

Kerry R. and Oliver MA, 2007. Comparing Sampling Needs for Variograms of Soil Properties Computed by the Method of Moments and Residual Maximum Likelihood. Geoderma 140, 383-396.

https://doi.org/10.1016/j.geoderma.2007.04.019

https://doi.org/10.1016/j.geoderma.2007.04.019

Kerry R. and Oliver MA, 2007. The Effects of Underlying Asymmetry and Outliers in

data on the Residual Maximum Likelihood Variogram: A Comparison with the Method of Moments Variogram.

Lark RM, 2000. A comparison of some robust estimators of the variogram for use in soil survey. European Journal of soil science 51(1): 137-157.

https://doi.org/10.1046/j.1365-2389.2000.00280.x

Lark RM., Hamilton EM, Kaninga B, Maseka KK, Mutondo M, Sakala GM, Watts MJ, 2017. Planning spatial sampling of the soil from an uncertain reconnaissance variogram. Soil 3(4): 235-244.

https://doi.org/10.5194/soil-3-235-2017

Liao KW, Guo JJ, Fan JC, Huang CL, Chang SH, 2019. Estimation of soil depth using Bayesian maximum entropy method. Entropy 21(1): 69.

https://doi.org/10.3390/e21010069

PMid: 33266785 PMCid: PMC7514177

Lichstein JW, Simons TR, Shriner SA, Franzreb KE, 2002. Spatial autocorrelation and autoregressive models in ecology. Ecological Monographs 72

https://doi.org/10.1890/0012-9615 (2002)072[0445: SAAAMI] 2.0.CO; 2 (3):445-463.

Manikandan S, 2010. Data transformation, Journal of Pharmacology Pharmacotherapeutics 1(2). DOI: 10.4103/0976500X.72373

https://doi.org/10.4103/0976-500X.72373

PMid: 21350629 PMCid: PMC3043340

Matheron G, 1963. Principles of Geostatistics. Economic Geology

https://doi.org/10.2113/gsecongeo.58.8.1246

Meul M. and Van Meirvenne M, 2003. Kriging soil texture under different types of nonstationarity. Geoderma 112(3-4): 217-233.

https://doi.org/10.1016/S0016-7061 (02)00308-7

Neerchal NK, Lacayo H, Nussbaum BD, 2008. Is a Larger Sample Size Always Better? American Journal of Mathematical and Management Sciences 28 (3-4): 295-307.

https://doi.org/10.1080/01966324.2008.10737 730

Obaid AN. and Mohammed MJ, 2020. A comparison of topological kriging and area to point kriging for irregular district area in Iraq. J. Mech. Cont. Math. Sci. 15(4): 105-124.

https://doi.org/10.26782/jmcms.2020.04.0000 9

Orton TG. and Lark RM, 2007. Accounting for the uncertainty in the local mean in spatial prediction by Bayesian Maximum Entropy. Stoch Environ Res Risk Assess 21:773-784. DOI 10.1007/s00477-006-0089-7.

https://doi.org/10.1007/s00477-006-0089-7

Orton TG. and Lark RM, 2009. The Bayesian maximum entropy method for lognormal variables. Stoch Environ Res Risk Assess 23:319-328. DOI 10.1007/s00477-008-0217-7

https://doi.org/10.1007/s00477-008-0217-7

Osborne W, 2010. Improving your data transformations: Applying the Box-Cox transformation. Practical Assessment, Research Evaluation 15(12).

http://pareonline.net/getvn.asp? v=15n=12.

Serre ML, Bogaert P, Christakos G, 1998. Computational investigations of Bayesian maximum entropy spatiotemporal mapping. In 4th International Association of Mathematical Geology (1): 117-122. Naples, Italy: De Frede Editore.

Society of Economic Geologists 58(8): 1246-1266.

Webster R. and Oliver MA, 1992. Sample Adequately to Estimate Variograms of Soil Properties. Journal of Soil Science

43:177-192. http://dx.doi.org/10.1111/j.13652389.1 992.tb00128.x

https://doi.org/10.1111/j.1365- 2389.1992.tb00128.x

Webster R. and Oliver MA, 2001. Geostatistics for environmental scientists. Wiley, Chichester

Yang Y, Zhang C, Zhang R, 2016. BME prediction of continuous geographical properties using auxiliary variables. Stochastic environmental research and risk assessment 30(1): 9-26.

https://doi.org/10.1007/s00477-014-1005-1

Zimmerman DL. and Zimmerman MB, 1991. A comparison of spatial semivariogram estimators and corresponding ordinary kriging predictors. Technometrics 33(1): 77-91.

https://doi.org/10.1080/00401706.1991.10484 771