# MATH460:Geostatistics, Project Geostatistical Modelling of Spatial Distribution of Lindane Concentrations in the European Continent

Ioannis Elmatzoglou

April 2, 2006

### Abstract

$\gamma$-*Hexachlorocyclohexane* ($\gamma$-*HCH* or *lindane*) has been extensively used worldwide as an agricultural insecticide , with most of the European countries having banned or severely restricted its use since the mid-1990s. European estimates suggest that as much as *135,000* tonnes was applied over the period *1970-1996* ; with the major emissions originating from France, Spain and the Netherlands. It is included in the list of chemicals that are considered to undergo long-range atmospheric transport, exhibit high overall environmental persistence, based on measurement data documenting its ubiquity in the environment. 62 air samples were collected across Europe for 6 weeks in June and July 2002 and among the compounds quantified was lindane, detected in almost all of the samples. We are developing a geostatistical approach in order to investigate initially the mechanism that describes the spatial distribution of this chemical's concentrations and in addition to make predictions of these concentrations in 6677 unsampled locations across Europe, using a geostatistical approach.

## 1  Introduction

*Persistent Organic Pollutants* (POPs) are chemical substances that persist in the environment, bioaccumulate through the food web, and pose a risk of causing adverse effects to human health and the environment. With the evidence of long-range transport of these substances to regions where they have never been used or produced and the consequent threats they pose to the environment of the whole globe, the international community has now, at several occasions called for urgent global actions to reduce and eliminate releases of these chemicals.

*Lindane - g-HCH [hexachlorocyclohexane]*, is included in the government "Red List" of dangerous substances. It has been in use as a broad range insecticide for 50 years, long enough to build up a significant body of evidence on its toxic and environmental hazards. It has caused deaths and poisonings in humans and there is authoritative recognition of its long term health effects including carcinogenic effects. It is a serious environmental contaminant and as well as being directly toxic to wildlife. It bio-accumulates along food chains. Scientific and anecdotal evidence links lindane with serious health problems including aplastic anaemia, the birth disorders C.H.A.R.G.E. and breast cancer. Lindane contaminates drinking water sources, it is highly persistent and travels long distances via atmospheric and oceanic currents.

Its major emissions originate from France, Spain and the Netherlands. In addition to the direct application of lindane there may also have been diffuse emission sources such as from disused chemical factories or/and emissions from dumps (*Manz et al., 2001*) revolatilisation from soils and other lindane-treated surfaces, as well as significant amounts "imported" into the European environment from the West and East (*Prevedouros et al., 2004*).

## 1.1    Analysis and Scientific Interest

The random and, in particular, the structural behavior of the spatial phenomena are, for convenience, considered within a geostatistical framework. From the geostatistical point of view, the spatial distribution of lindane is viewed as a result of a stochastic process characterized by some geostatistical parameters, the most important of which synthesizes the continuity structure derived from the data. Once an adequate number of parameters of stochastic process is known, and some hypothesis have been formulated, this probabilistic framework allows us to analyze and quantify the uncertainty characterizing the spatial distribution of the studied variable. In this work we have adopted a geostatistical procedure to study the spatial distribution of lindane's concentrations in the European continent using a sample of 62 measurements.

Estimates of lindane's air concentrations, is something of vital importance for the decision making process and the future potential (further) restriction in its emission levels, in the countries in which it is not already banned. Interest also lies in understanding the complex array of factors controlling its concentrations. For this reason we are going to investigate whether these quantities can be influenced by temperature and precipitation in each particular location.

## 1.2    Background on Measurement Datasets

**Concentrations:** A passive air sampling campaign (PASAE) was carried out across Europe for 6 weeks in June and July 2002. A total of 62 polyurethane foam (PUF) discs were deployed in urban and rural sites covering Iceland and Portugal in the west to Russia, Cyprus and Kazakhstan in the east (*Jaward et al., 2004*). Among the compounds quantified was lindane, detected in almost all of the samples. This was the first time that concurrently sampled ambient air data of such spatial scale have been reported.

**Temperature and Precipitation:** Monthly averaged global climate data from the International Institute for Applied System Analyses (*IIASA; Laxenburg, Austria*) which includes precipitation and temperature data were used (*Leemans and Cramer, 1991*). The dataset uses standardized climate records from up to eight different sources, interpolated and smoothed to fit a one-half degree latitude/longitude terrestrial grid surface.

# 2    Exploratory Data Analysis

## 2.1    Viewing the Data

In this first stage of the exploratory analysis of the data, we make a plot of the lindane measurements (response variable) in relation to their locations (*figure 1*).
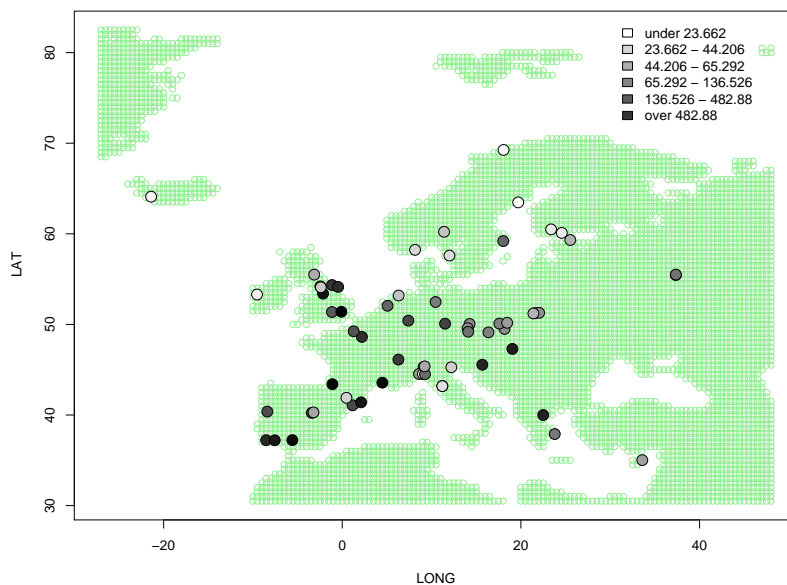
**Figure 1**: *Concentrations of g-HCH in the sample locations*

We are able to see that the greatest measurements of the chemical are on the south-west Europe especially in the regions covering France, Spain, Portugal, as well as some of the central European countries. However, it is not possible for us to make any assessment of the spatial correlation or to distinct any non-spatial trends. In the following plot we just visualize the annual precipitation data for the year 2002 and the sampling locations.
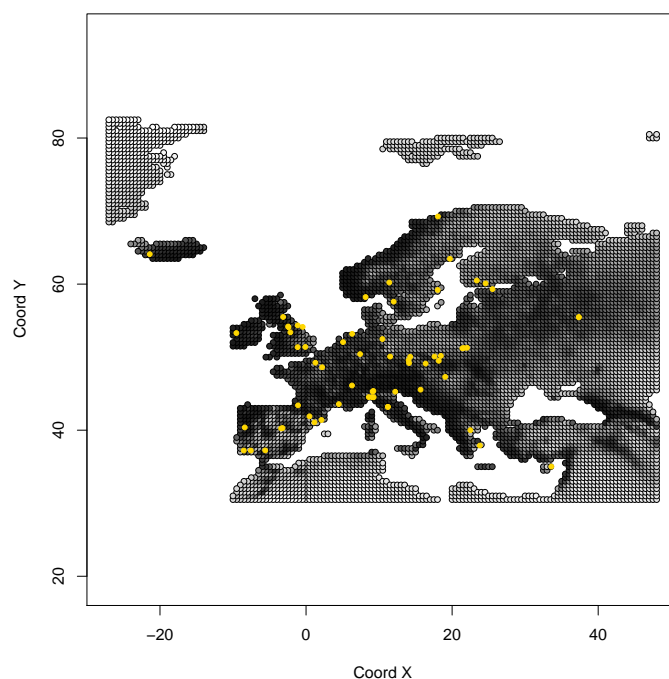


**Figure 2**: *Precipitation in Europe and sample points of Lindane (g-HCH)*

We can observe that the greatest precipitation levels for this year took place basically in the central and southern Europe. It is impossible to distinct any link between precipitation and lindane levels. In the next part we are focusing on finding an appropriate model for our data.

## 2.2   Analysis

Gaussian stochastic processes are widely used in practice as models for geostatistical data. They are used as convenient empirical models which can capture the a wide range of spatial behavior according to the specification of the correlation structure. One very good reason for concentrating on the gaussian models is that they are quite convenient and uniquely tractable as models for dependent data. So, our analysis will be performed under the gaussian assumption for the distribution of our response variable. By denoting $\mathbf{Y}$ to be the vector of these measurements:

$$Y \sim MVN(\mu, \mathbf{V})$$

, where $\mu$ is the vector of the *62* means and $\mathbf{V}$ is the $(62 \times 62)$ variance-covariance matrix which equals to $\sigma^2 R + \tau^2 I$ and represents the dependence between the measurements. Specifically, $\tau^2$ expresses the mean square errors between the actual concentrations of lindane and our measurements and $R$ the $(62 \times 62)$ correlation matrix with elements $r_{ij} = \rho(||\omega_i - \omega_j||)$., where $\omega_i$'s are the sample locations of the chemical *i=1,2, . . .,62* in the 2-dimensional space.

But how realistic is the normality assumption in our particular case? All we have to do is to make a histogram of our *62* data. In the part (A) of the next figure a histogram of the data is provided:
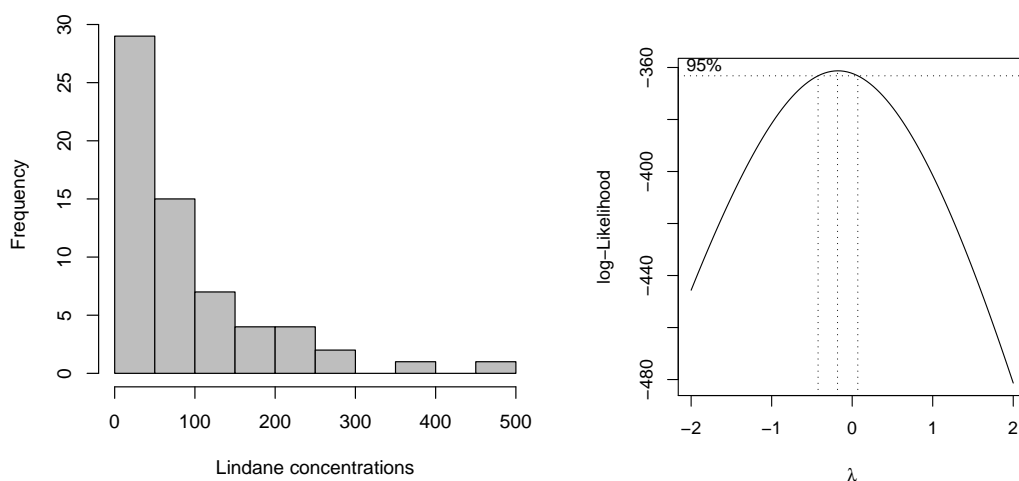


**Figure 3, A & B**: *A:Histogram of the chemical, B:Profile Likelihood for the parameter $\lambda$ of Box-Cox transformation*

It is obvious that the data don't look like Gaussian at all and instead, it seems that they were "generated" by an asymmetric distribution. Nevertheless, the range of applicability of the Gaussian model can be extended by assuming that the model holds after a marginal transformation of the response variable. *Box* and *Cox* have proposed the following parametric family of transformations (*Box & Cox 1964*):

$$Y^* = \begin{cases} \frac{Y^\lambda - 1}{\lambda} & : \lambda \neq 0 \\ log(Y) & : \lambda = 0 \end{cases}$$

, where a particular choice of $\lambda$ can lead to an empirical Gaussian approximation. The best approximation can be found by maximizing the *profile likelihood* of $\lambda$. We maximized this likelihood after taking into-consideration the average annual temperature and precipitation of the year 2002, as well as their average values for the particular month (July) where the sampling took place. This profile likelihood can be viewed in the second part of the *figure 3*. As an estimate of $\lambda$ we chose the value *-0.15*. What we are going to do is to fix the parameter $\lambda$ to be equal to this estimate. Then we are going to transform the data and after making inferences on the transformed scale, we are going to re-transform them into the original scale (*Christensen, Diggle, Ribeiro 2000*). We can get a view of the 'new' data through these following plots.
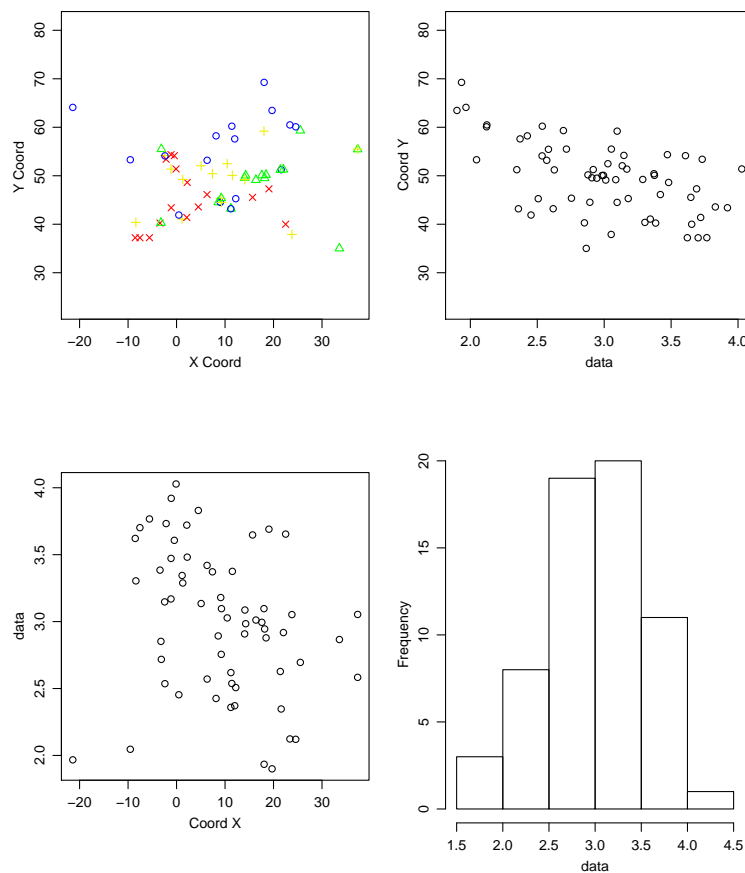


**Figure 4**: *Box-Cox transformation of the data, using $\lambda = -0.15$*

The Gaussian approximation seems to be quite good, which means that we can work under the Gaussian assumptions. Moreover, we can observe in the second scatterplot a small evidence of spatial trend and particularly a north-south trend with higher responses concentrated towards the southern parts of the European continent.

A vital part of the geostatistical modelling and exploratory analysis of the data consists of investigating the correlation structure of the data and making appropriate assumptions for this structure. The first plots were not so helpful in this assessment. In the next part of the analysis we are dedicated in specifying the most suitable form of the correlation in our data.

## 2.3 Correlation Structure

A good visual assessment of the spatial correlation of our response variable can be done by means of the empirical variogram. The empirical variogram is a plot consisting of points generated by the following theoretical function:

$$V(x, x + h) = \frac{1}{2} Var\{Y(x) - Y(x + h)\}$$

, that is the variance of difference in the response variable measured in two locations that are h distance apart. This is directly linked with the correlation structure of the process and particularly in the case of a stationary process this equals to:

$$V(h) = \frac{1}{2} E([Y(x) - Y(x + h)]^2) = \tau^2 + \sigma^2 \{1 - \rho(h)\}$$

, where $\tau^2$ reflects the variance in the measurement errors, $\sigma^2$ the variance of the stochastic process and $\rho(h)$ the correlation between two given points that are $h$ distance apart. An estimation of this theoretical line can be performed by calculating:

$$\frac{1}{2N(h)} \sum_{i=1}^{N(h)} \left[ Y(x) - Y(x + h) \right]^2$$

, for many values of $h$, where $N(h)$ is the number of pairs of measurements that are h distance apart in the 2-dimensional space. Note that in this empirical function $h$ denotes usually a range of values. After calculating this empirical line we can have a quite clear "image" for a correlation pattern in our sample which is an estimate of the real correlation structure.

In the case of non-stationarity in our process, the same analysis can be applied but after detrending the data. For this reason we tried to fit a kind of trend surface model by considering all the available covariates. Then we tried to detect the covariance structure of the unobserved stochastic process $S$ by applying the same things described before but now working with the residuals of this model, which has as covariates the most significant explanatory variables found:

$$r_i = Y_i - \hat{\mu}(x_i)$$

$$v_{ij} = \frac{(r_i - r_j)^2}{2}$$

The empirical variogram can be estimated by calculating the average of $v_{ij}$'s for a given range of distance e.g. the average of $v_{ij}$'s that correspond to pairs of measurements that are h distance apart, where $h_k \leq h \leq h_k + \ell$ , and $\ell$ denotes the length of each of the k defined bins:

$$\gamma(h) = \sum_{z=1}^{N(h)} \frac{v_{ij}}{N(h)}$$

Note that $N(h)$ is the number of pairs of residuals that correspond to measurements that are $h$ distance apart in the 2-dimensional space.

In our case, the most significant covariates were found to be average temperature (year 2002) , while of secondary importance were the average precipitation of the same year as well as the average temperature in month July, where the sampling took place. The variogram produced by the residuals of this *ols* model can be seen in the following plots.
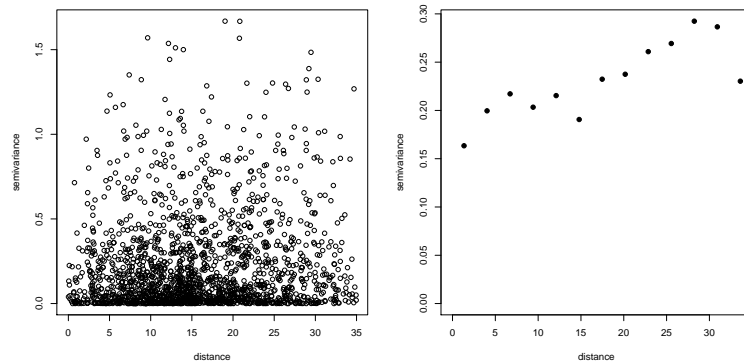
**Figure 5**: *Empirical variogram of the transformed and detrended data*

In the first figure we are viewing the $\frac{(r_i - r_j)^2}{2}$ points plotted against their distance ($h_{ij} = ||x_i - x_j||$), and in the second plot there is the average of these points. This enables us to have a clearer image about the kind of correlations. There are various families that describe this correlation structure of the process. However, we are going to focus on three types of potential correlation and trying to figure out which is the most appropriate for our case. These are the *exponential*, *spherical* and the *gaussian*, which are given by the following relations:

Exponential

$$\rho(h) = e^{-\frac{h}{\phi}}, \phi > 0$$

Gaussian

$$\rho(h) = e^{-\left(\frac{h}{\phi}\right)^2}, \phi > 0$$

Spherical

$$\rho(h) = \left\{ \begin{array}{ll} 1 - \frac{3}{2}(h/\phi) + \frac{1}{2}(h/\phi)^3 & : 0 \leq h \leq \phi \\ 0 & : h > \phi \end{array} \right.$$

What follows is to fit a model with these three kinds of spatial covariance structure, compare the results and deciding which is the most suitable one.

In geostatistics there are three different ways of making inferences about the unknown parameters. We can follow either the classical or the bayesian approach but another way is to use the variogram as an inferential tool. Here we are focusing on the classical inferential methodology.

## 3   Model Fitting

We are adopting the model described in *section 2.2* and we remind that as a response variable we now have the transformed (*according to the Cox-Box*) observed values of the chemical. By taking into account the covariates influencing the mean of the unobserved stochastic process $S$, then we have the following Gaussian model with linear specification trend:

$$Y^* \sim MVN(D(\omega)\beta^T, \sigma^2 R(\phi) + \tau^2 I)$$

, where $Y^* = \{Y_1^*, Y_2^*, \ldots, Y_{62}^*\}$, $D(\omega)$ is the $(6677 \times 3)$ covariates matrix and $\beta = \{\beta_0, \beta_1, \beta_2\}$, the parameters that correspond to each of the *3* coefficients: the intercept and the two

covariates which are the average annual precipitation and temperature. Note that $Y^* = \frac{Y^{(-0.15)}-1}{(-0.15)}$.

The log-likelihood function can be written as:

$$L(\beta, \tau^2, \sigma^2, \phi) = -0.5\Big\{log(2\pi) + log\{|(\sigma^2 R(\phi) + \tau^2 I)|\} + (y^* - D(\omega)\beta^T)^T(\sigma^2 R(\phi) + \tau^2 I)^{-1}(y^* - D(\omega)\beta^T))\Big\}$$

We maximized this function for the three different types of spatial correlation structure described before. The values of the parameters that maximize it in each particular case, are given by the following table:

| | $\beta_0$ | $\beta_1$ | $\beta_2$ | $\tau^2$ | $\sigma^2$ | $\phi$ | log-likelihood |
|---|---|---|---|---|---|---|---|
| Exponential | 2.3177 | -0.0016 | 0.0698 | 0.0826 | 0.1694 | 2.7569 | -321.3 |
| Spherical | 2.2991 | -0.0013 | 0.0754 | 0.0726 | 0.1624 | 7.8839 | -321.1 |
| Gaussian | 2.2501 | -0.0008 | 0.0737 | 0.0898 | 0.1649 | 4.5170 | -320.9 |

**Table 1** *Maximum likelihood estimates for the three different models in respect to the correlation Structure*

The estimation of the vector of $\beta$ parameters seems to be quite similar between the *3* cases, as well as the estimation of the $\sigma^2$ (*sill*) covariance parameter. This does not hold for the case of the *nugget* ($\tau^2$) and especially for the *range* covariance parameter $\phi$. After performing some Likelihood Ratio Tests, annual average temperature was found to be the most significant explanatory variable. Additionally, we estimated the parameters of the likelihood function considering also the fact that we may have anisotropy. After performing a *LRT*, and under the hypothesis of not having anisotropy, the $\chi^2$ test for *2* degrees of freedom (for the two additional parameters in the model) gave a relatively big *p-value*, which led us accept the null hypothesis of non-anisotropy in our data.

The table give us information about the value of the *log-Likelihood* for each of the three cases. Its biggest value corresponds to the case where we have a *Gaussian* spatial correlation pattern, but the difference in these values is not so big. An alternative way of exploring the appropriateness of each of the three different correlation schemes, can be done by comparing the empirical variogram with the *3* corresponding fitted lines produced by substituting the maximum likelihood estimates to the *3* theoretical function forms:
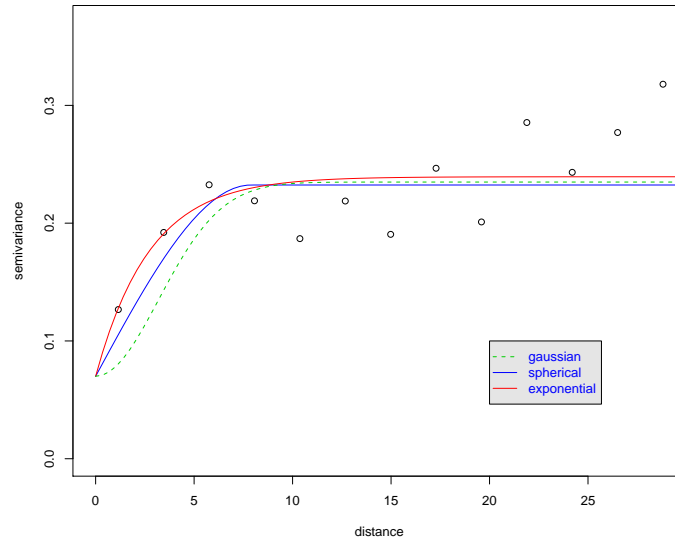
**Figure 6**: *Comparing the empirical variograms with 3 fitting variogram functions with parameters obtained by ML estimation that correspond to different variogram structures*

The main differences between these lines is at the first few values of distance. Obviously the difference in the *3* specified correlation functions and the variation in the *ML* estimations of the covariance parameters, can be reflected in this plot. Although the *exponential* function is the one that gives the lowest log-likelihood value for our model, is the one that seems to "explain" the spatial correlation in the most appropriate way. And this is the reason why we are going to use it in order to make predictions over the studied area.

# 4   Spatial Prediction

The model found suitable in terms of interpreting the stochastic behavior of our data, was:

$$Y^* \sim MVN(D(\omega) \cdot \hat{\beta}^T, \hat{\sigma}^2 R(\hat{\phi}) + \hat{\tau}^2 I)$$

, where $\hat{\beta} = \{2.3177, 0.0698\}$, the parameters that correspond to the intercept and to the average temperature in each particular point,$\hat{\tau}^2 = 0.0826$, $\hat{\sigma}^2 = 0.1694$, $\hat{\phi} = 2.7569$ , $R$ is the $(62 \times 62)$ matrix with elements $r_{ij} = exp(-\frac{|\omega_i - \omega_j|}{\hat{\phi}})$ , where $\omega = \{\omega_1, \omega_2, \ldots \omega_{62}\}$ $(\omega_i \subset \Re^2)$.

That means that if our assumptions are true and our estimations valid, then the lindane's concentrations can be viewed as a the result of discipline to the following stochastic mechanism:

$$S \sim SGP(D(\mathbf{x})\hat{\beta}^T, \hat{\sigma}^2)$$

, where $\mathbf{x} \subseteq \Re^2$.

If this is the case, then for a vector of prediction points: $\mathbf{x} = (x_1, x_2, \ldots, x_n)$, it holds that:

$$S(\mathbf{x}) \sim MVN(D(\mathbf{x})\beta^T, \sigma^2)$$

This enables us to find an expression about the joint density function of the 2 random vectors $Y^*$ and $S(\mathbf{x})$, which is:

$$\begin{bmatrix} S(\mathbf{x}) \\ Y^* \end{bmatrix} \sim MVN \begin{bmatrix} \mu_1(\mathbf{x}) & Var(S(\mathbf{x})) & Cov(S(\mathbf{x}), Y^*) \\ \mu_2(\omega) & , & Cov(Y^*, S(\mathbf{x})) & Var(Y^*) \end{bmatrix}$$

, with $\mu_1(\mathbf{x}) = D(\mathbf{x})\hat{\beta}^T$, $\mu_2(\omega) = D(\omega)\hat{\beta}^T$, $Var(S(\mathbf{x})) = \hat{\sigma}^2 R(\hat{\phi})$, $Var(Y^*) = \hat{\sigma}^2 R(\hat{\phi}) + \hat{\tau}^2 I$, $Cov(S(\mathbf{x}), Y^*) = \hat{\sigma}^2 r'$ and $Cov(Y^*, S(\mathbf{x})) = \hat{\sigma}^2 r$.

Accordingly, in the case that $\mathbf{x}$ corresponds to a single unsampled point, $x_i$, for $S(x_i)$, $i = 1, 2, \ldots, n$, it holds that the distribution of the chemical's concentrations in this particular point is given by:

$$S(x_i)|Y^* \sim Normal(M_i, \Sigma_i)$$

, with first and second moments:

$$M_i = D(x_i)\hat{\beta}^T + \hat{\sigma}^2 \mathbf{r'}(\hat{\sigma}^2 R(\hat{\phi}) + \hat{\tau}^2 I)^{-1}(Y^* - D(\omega)\hat{\beta}^T)$$

$$\Sigma_i = \hat{\sigma}^2 - \hat{\sigma}^2 \mathbf{r'}(\hat{\sigma}^2 R(\hat{\phi}) + \hat{\tau}^2 I)^{-1}\sigma^2 \mathbf{r}$$

, where $\mathbf{r}$ is the vector of the elements $r_j = exp\left(-\frac{|x_i - \omega_j|}{\hat{\phi}}\right)$ for $j = 1, 2, \ldots, 62$. This particular property of the Gaussian distribution provides us with an easy and convenient way of doing predictions in the locations where the value of the response variable is unknown. This way of doing predictions is called *simple kriging*. Furthermore, it provides us with a degree of quantification of the uncertainty associated with estimates. So, in our particular case, we are going to make predictions to 6677 unsampled points by conditioning on the data already observed:

$$E(S(x_i)|Y^*) = M_i$$

$$\text{Var}(S(x_i)|Y^*) = \Sigma_i$$

, for $i = 1, 2, \ldots, 6677$

We can see that kriging makes maximum use of the information available from known values. The *6677* predicting data were back-transformed into their original scale and the results of the predictions can be viewed into the following two figures:
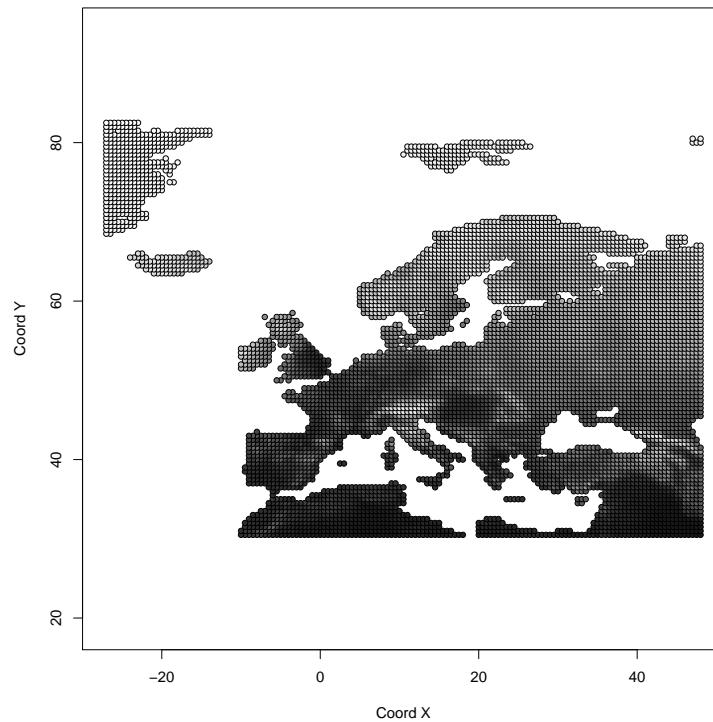
**Figure 7**: *Spatial predictions-Simple kriging with trend*

, where darker values represent high concentrations of the chemical. We can have a clearer image of the predictive values by plotting the 20%, 40%, 60% and 80% quantiles of the predictive values with different color:
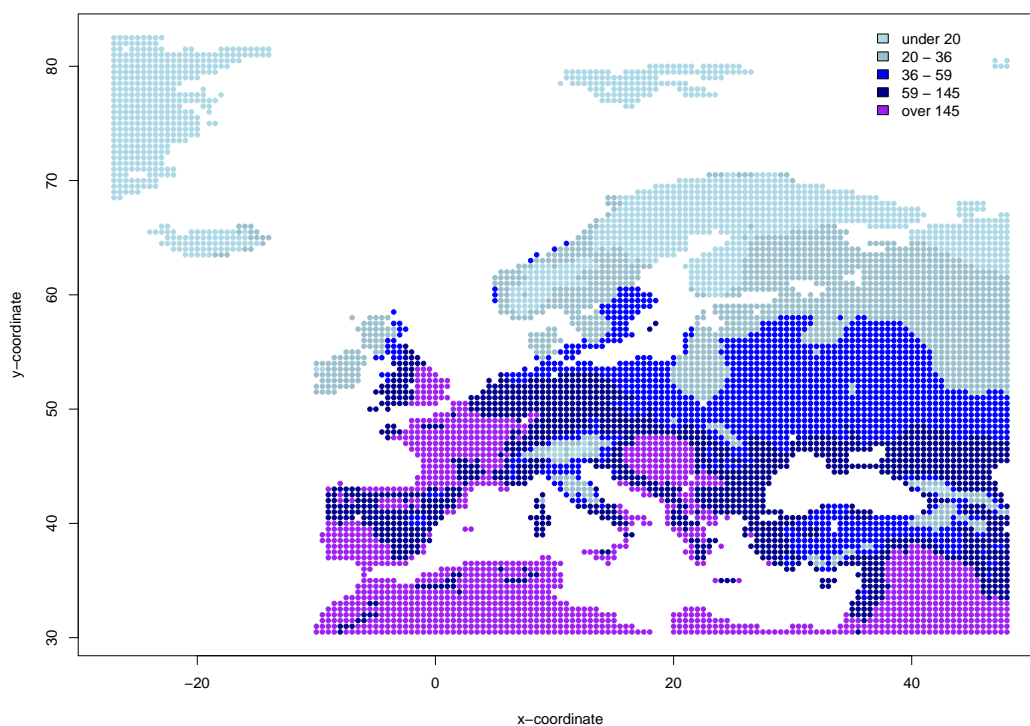
**Figure 8**: *20%, 40%, 60%, 80% quantiles of the predictions*

So, by making maximum use of the available information, we can say that the highest expected concentrations of lindane can be found basically in the areas covering Spain (and a part of Portugal), France, Belgium , Netherlands, south-east England, Hungary, Slovenia, Servia, Romania, the northern part of Greece and the southern part of Italy. High concentrations of the chemical appear also to be in the north-African countries (*see Discus.*). The lowest expected concentrations of lindane are found mainly in the North European and Scandinavian countries such as Norway, Finland and Sweden, Iceland as well as the northern part of Italy. The variance quantiles plot was found very similar, which means that our uncertainty increases as we are moving towards the southern parts of the studied area.

# 5 Results-Conclusion

With the main interest lying in investigating the stochastic spatial behavior of the persistent organic pollutant lindane ($\gamma$-$HCH$), we were based on a random sample of 62 concurrent observations around Europe. From the beginning it was obvious that the data were not normally distributed. Nevertheless, we transformed them by making use of the Box-Cox transformation and we got a good Gaussian approximation. So, working under the Gaussian assumptions, we tried initially to detect the correlation scheme of the unobserved stochastic process by inspecting the way that our measurements were correlated. For this reason we constructed the empirical variogram of the data and examined the suitability of three different kinds of correlation: exponential, spherical and Gaussian. The behavior of the estimated variogram at the first lags led us preferring the exponential one. By following the classical inferential approach, we focused on estimating the parameters of the model that we specified. The latter was a multivariate Gaussian with linear specification trend model with covariates the average temperature and precipitation in the locations (and at the year) where the sampling took place. Average temperature was found only to affect the mean of the stochastic process. By being based on these estimations, we made predictions in *6677* locations by means of the simple kriging method. The predictions were re-transformed back to their original scale and we viewed the results on the figures. As expected, the highest predictive values in the concentrations of the chemical were found to be in countries consisting at highest percentage by agricultural regions, such as France, Spain and as well as other central and southern European Countries, where lindane has been extensively used as an insecticide. The lowest expected concentrations of the chemical were basically in the north European and Scandinavian countries, but also in some other smaller southern regions, such as the north part of Italy.

# 6 Discussion

We can see in the quantiles map a downwards trend in the predictive values, while we see the north-African countries belonging to the areas with the highest predictions. There are two basic explanations for this fact. First, we don't have any measurements in these areas, which means that as we move further and further from our observations, the expected values tend to be equal to the mean of the specified stochastic process and its increasing downwards trend, a part of which, is caused by the increase in the temperature (which is relatively more increased in the southern countries). Secondly, as this trend was found by taking into

consideration measurements not belonging to the northern African countries, we are unsure about its existence if we move further to the southern direction. A plausible example is that if this trend is actually caused by the existence of more agricultural regions in the southern European countries (where the pesticides are more intensively used), it is almost impossible for us to expect that the concentrations in the north-African countries will still be high (where there are not such kind of regions). If this is the case, then the information corresponding to the areas very far away from our observations can be quite misleading as we haven't taken into account such kind of spatial factors in our model. In the following figure there is a contour plot with the kriging predictions produced under the assumption of constant mean in our process, but by letting the covariance parameters to be the same as in the linear trend case.
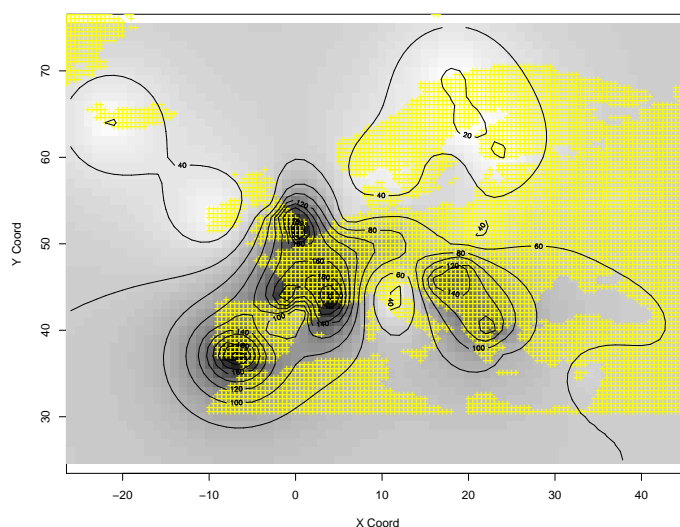


**Figure 9**: *Spatial predictions-Simple kriging without trend*

Some "bad" results caused by a specification of a downwards trend (increasing temperatures to the south countries) are now disappeared. It is believed that this plot approximates in a better way the reality in the southern countries of the map. And that because there is probably an inverse trend, that is when we are moving from the very southern areas of our region of study to the more central ones. Not only the percentages of agricultural regions increases but also some other unconsidered factors such as the increase in the number of the factories, which is a secondary source of lindane's emissions. By not considering these opposite spatial correlated factors (downwards (i.e temperature) and upwards factors), we let them be incorporated into the stochastic process and the results of predictions are probably better for the regions that are not very close to our sample locations.

Something that could explain a part of the non-spatial variability in the lindane's concentrations could have also been an index of the restriction levels in each country.

Finally, another problem that arises is that of multicollinearity. Temperature increases towards the smallest values of the *y-axis* and the fact that it was found to be significant may be caused by the potential case that it is correlated to other significant covariates that probably increase towards the same direction. So the effect of temperature in this chemical concentrations raises many doubts.

# References

[1] Diggle, P.J., Ribeiro Jr., P.J. (2006) Model-Based Geostatistics *(Draft Version)*

[2] Jaward F.M., Farrar N.F.,Harner T., Sweetman A.J. and Jones K.C. Passive Air Sampling of PCBs, PBDEs, and Organochorine Pesticides Accross Europe, *Enviromental Science and Technology 2004, 38, 34-41*

[3] Prevedouros K., MacLeod M., Jones K.J., Sweetman A.J., Modelling the Fate of Persistent Organic Pollutants in Europe: Parameterisation of a Gridded Distribution Model, *Environmental Pollution 128 (2004) 251-261*

[4] Fabbri P., Trevisani S., Spatial Distribution of Temperature in the Low-Temperature Geothermal Euganean Field (NE Italy): A Simulated Annealing Approach, *Geothermics 34(2005) 617-631*

[5] Ribeiro Jr., P.J. , Diggle, P.J. (2001) geoR: A package for geostatistical analysis. *R-NEWS, Vol 1, No 2, 15-18. ISSN 1609-3631.*