



<https://www.menti.com/alz6g6hw4mr7>

- In the last two weeks, we looked at time series - data which vary over *time*.
- We noted that observations closer together in time were likely to be correlated, and we discussed methods for accounting for that correlation.
- Over the next two weeks, we will look at data which vary over *space*.
- This takes us into the field of spatial statistics, which focuses on developing methods to account for spatial correlation.

Introduction to Spatial Statistics



<https://www.menti.com/alz6g6hw4mr7>

- Spatial data are data which have any form of geographical information attached to them.
- Spatial data are becoming increasingly frequent as a result of new technology (eg satellite weather data, GPS enabled devices including phones, CCTV systems).
- Most environmental variables of interest will vary over space to some degree, for example global temperature, air pollution, species distribution.
- We will typically use this spatial information to help us understand the relationship between our data points

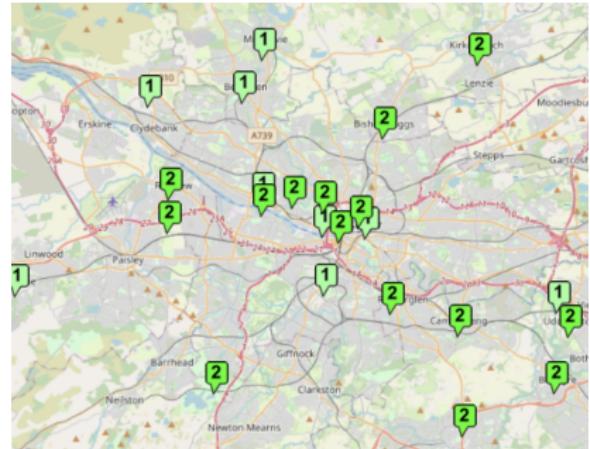
- Generally, each datapoint will contain a measured value in addition to information about their location (typically represented by some form of coordinate system).
- However, in some cases the locations themselves will be our data (eg the positions of trees in a forest).
- We typically assume that we have some degree of spatial autocorrelation present in our data.
- Tobler's first law of geography states that:
“Everything is related to everything else, but near things are more related than distant things”

- Air pollution levels are measured at monitoring stations around Glasgow.
- 1 is “*Low*” and 10 is “*Very High*”.

17th July 2022



17th January 2023



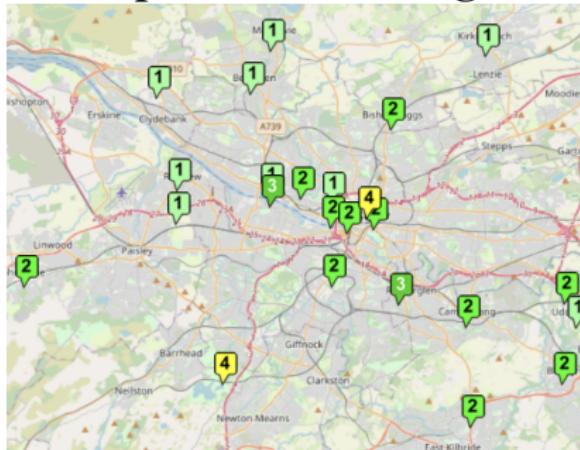
<https://www.scottishairquality.scot>

- The overall goal of any piece of spatial analysis is to understand the spatial patterns in our data.
- This could involve any of the following:
 - Estimating differences in mean, variance or some other summary statistic over space.
 - Predicting the value at some unobserved location.
 - Identifying “hotspots” with a high (or low) value compared to the rest of the region.
- Typically this involves understanding and accounting for spatial autocorrelation in our data.

- There are three main categories of spatial data which we will look at during this course.
- These are geostatistical data, point processes and areal (or lattice) data.
- We will briefly introduce each of them and identify the main topics where they are used for environmental data.
- Note that there is a separate Spatial Statistics course which covers each of these types of data in more detail.

- Measurements are taken at a set of fixed locations to measure some continuous process.

Air pollution in Glasgow



Bathing water quality in Mallorca



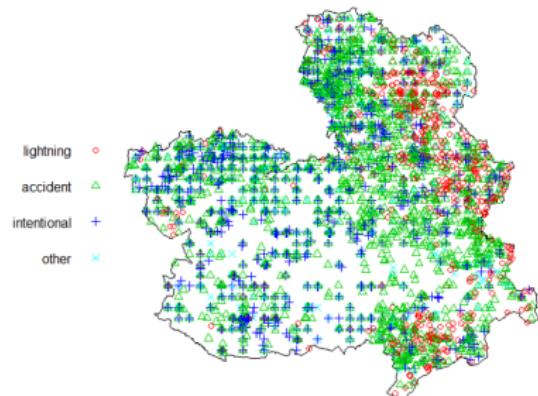
<https://www.eea.europa.eu/themes/water/interactive/bathing/state-of-bathing-waters>

- We measure the locations where events occur (eg trees in a forest, earthquakes) and the coordinates are our data.

Locations of flowers

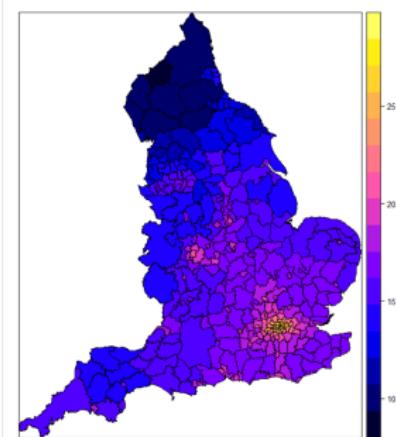


Forest fires in Castilla La Mancha, Spain

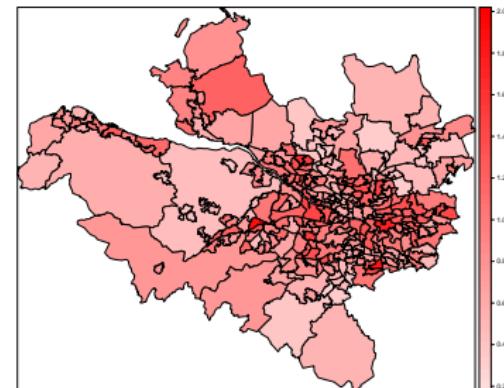


- Our data are measured as summaries across a set of discrete, non-overlapping spatial units (eg postcode areas, council regions).

Pollution in English local authorities



Human respiratory disease in Glasgow



Introduction to Geostatistical Data

- Geostatistical data are the most common form of spatial data found in environmental setting.
- We regularly take measurements of an environmental variable of interest at a set of fixed locations.
- This could be data from samples taken across a region (eg water depth in a lake) or from monitoring stations as part of a network (eg air pollution).
- In each of these cases, our goal is to estimate the value of our variable across the entire space.

- Let D be our two-dimensional region of interest.
- In principle, there are infinite locations within D , each of which can be represented by mathematical coordinates (eg latitude and longitude).
- We can identify any individual location as $s_i = (x_i, y_i)$, where x_i and y_i are their coordinates.
- We can treat our variable of interest as a random variable, Z which can be observed at any location as $Z(s_i)$.

- Our geostatistical process can therefore be written as:

$$\{Z(s); s \in D\}$$

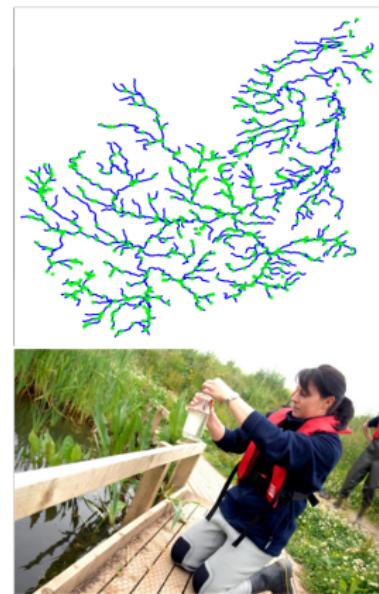
- In practice, our data are observed at a finite number of locations, m , and can be denoted as:

$$z = \{z(s_1), \dots, z(s_m)\}$$

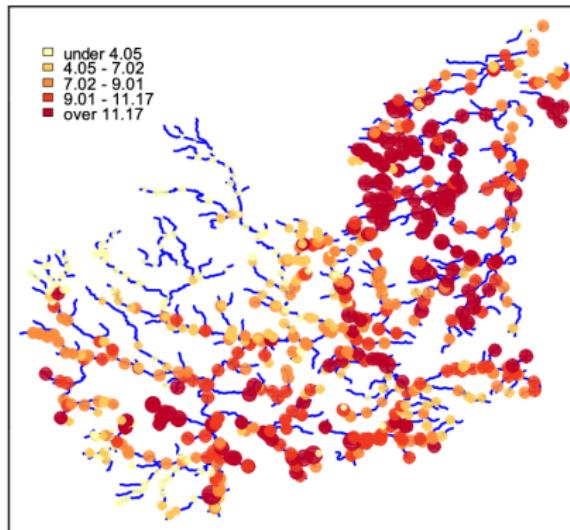
- We have observed our data at m locations, but often want to predict this process at a set of unknown locations.
- For example, what is the value of $z(s_0)$, where s_0 is an unobserved site?

- There are two main steps in classical geostatistical analysis.
 - 1 How do I produce a statistical model for the data?
 - 2 How do I use my model to estimate quantities of interest?
- The first part requires us to think about how our measured datapoints relate to each other - in other words, to understand spatial autocorrelation.
- The second part requires us to use that information to predict the value at unmeasured locations, and then to produce maps or summary statistics based on this.

- We are interested in nitrate levels in the River Trent.
- A set of locations in the river network are sampled.
- There are UK/EU directives for safe nitrate levels in water bodies such as rivers.
- Our goal is to answer a **policy** question - which areas breach 50mg/l limit?
- This requires us to estimate the levels at unmeasured locations.



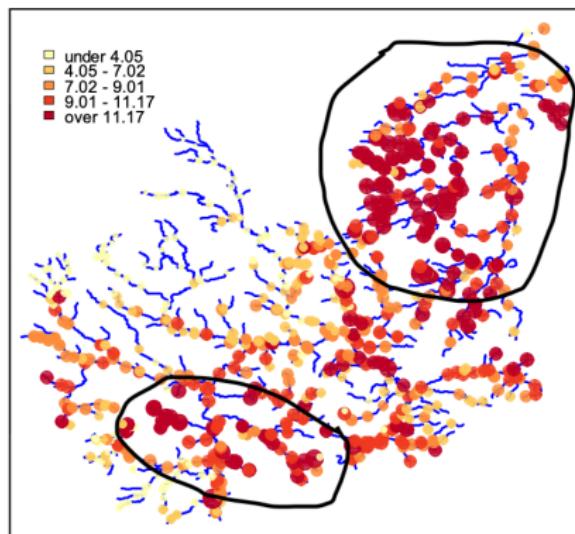
- The plot below shows the average nitrate levels from 2003-2010 at each of our observed locations.
- Darker colours and larger plotting points correspond to higher levels of nitrate.





<https://www.menti.com/alz6g6hw4mr7>

- Visually, it appears that the highest levels are located in the north-east and far south of the region.
- However, we require a statistical model to measure this objectively.



- The key challenge in modelling geostatistical data is understanding **correlation**.
- Typically observations close together in space will be more similar than those which are further apart.
- Spatial correlation is usually driven by some unmeasured confounding variable(s) - for example, air pollution is spatially correlated because nearby areas tend to experience similar traffic levels.
- It is important that we account for these correlations in our analysis - failing to do so will lead to poor inference.

- For a set of geostatistical data $\mathbf{z} = \{z(\mathbf{s}_1), \dots, z(\mathbf{s}_m)\}$, we can consider the general model:

$$Z(\mathbf{s}_i) = \mu(\mathbf{s}_i) + e(\mathbf{s}_i)$$

- Here $\mu(\mathbf{s}_i)$ is a mean function which models trend and covariate effects.
- Then $e(\mathbf{s}_i)$ is the error process which accounts for any spatial correlation which exists after accounting for $\mu(\mathbf{s}_i)$
- Spatial statistics is therefore often focused on understanding the process for $e(\mathbf{s}_i)$.

- We have observations at m locations

$$\mathbf{z} = \{z(\mathbf{s}_1), \dots, z(\mathbf{s}_m)\}.$$

- We want to use these to obtain an estimate of $Z(\mathbf{s}_0)$ where \mathbf{s}_0 is an unobserved location.
- How do we model the spatial dependence between our observed sites $\mathbf{s}_1, \dots, \mathbf{s}_m$?
- What does this tell us about the dependence between our observed sites and our unobserved site \mathbf{s}_0 ?

Variograms

- Spatial dependence is commonly modelled by a function known as a **variogram**.
- The variogram is similar in many ways to the autocorrelation function used in time series modelling.
- In simple terms, it is a function which measures the difference in the spatial process between a pair of locations a fixed distance apart.
- In order to define the variogram, it is important to first understand some more features of a geostatistical process.

- If we have a geostatistical process $\{Z(s); s \in D\}$, then its mean can be expressed as

$$\mu_z(s) = E[Z(s)] \text{ for all } s \in D.$$

- Our process Z is then said to be Gaussian if our random variable at the set of observed locations is multivariate normal.
- In other words, Z is Gaussian if

$$\{Z(s_1), \dots, Z(s_m)\} \sim \text{MVN}(\mu_z(s), C_z(s)).$$

- Similarly, the covariance can be expressed as:

$$\begin{aligned}C_z(s, t) &= \text{Cov}(Z(s), Z(t)) \\&= E[(Z(s) - \mu_z(s))(Z(t) - \mu_z(t))]\end{aligned}$$

- Here, the covariance measures the strength of the linear dependence between $Z(s)$ and $Z(t)$.
- As usual, we can compute the variance of $Z(s)$ as a special case of the covariance where $s = t$.

- Our geostatistical process can be described as **weakly stationary** if the following criteria are met:
 - 1 $E[Z(s)] = \mu_z(s) = \mu_z$ - a finite constant which does not depend on s .
 - 2 $C_z(s, s + h) = C_z(h)$ - a finite constant which can depend on h but not s .
- Condition 1 states that our mean function must be constant in space, with no overall spatial trend.
- Condition 2 states that for any two locations, their covariance depends only on how far apart they are (their **spatial lag**, h), not their absolute position.

- A geostatistical process is said to be **isotropic** if the covariance function is *directionally invariant*.
- This means that the covariance between two points a distance h apart is the same no matter which direction you travel in.
- Mathematically, this can be denoted by

$$C_z(\mathbf{h}) = C_z(||\mathbf{h}||).$$



<https://www.menti.com/alws53ba6bzm>

- In this course, we will only look at Gaussian, weakly stationary and isotropic processes.
- This means that:
 - This means our random variables across our observed locations follow a multivariate Gaussian distribution.
 - The mean of this distribution is constant over space.
 - The covariance of this distribution is only a function of the lag, not the position or direction of the points.
- Other models do exist for more complex processes, but we will not explore these.

- The function describing the dependence between values of our process Z separated by different lags is known as the **autocovariance function**.
- This is similar to the autocorrelation function (ACF) used for temporal data.
- In geostatistical models, a variant of this known as a **variogram** is typically used to estimate spatial relationships.

- The variogram measures the variance of the difference in the process at two spatial locations s and $s + \mathbf{h}$.
- The variogram is defined as

$$\text{Var}[Z(s) - Z(s + \mathbf{h})] = E[(Z(s) - Z(s + \mathbf{h}))^2] = 2\gamma_z(\mathbf{h}).$$

- Here, $2\gamma_z(\mathbf{h})$ is the variogram, but in practice we use the **semi-variogram**, $\gamma_z(\mathbf{h})$.
- We use the semi-variogram because our points come in pairs, and the semi-variance is equivalent to the variance per point at a given lag.

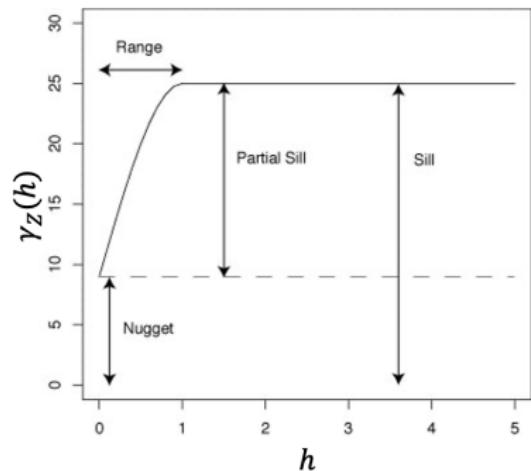
- When the variance of the difference $Z(s) - Z(t)$ is relatively small, then $Z(s)$ and $Z(t)$ are similar (spatially correlated).
- When the variance of the difference $Z(s) - Z(t)$ is relatively large, then $Z(s)$ and $Z(t)$ are less similar (closer to independence).
- If our process is *weakly stationary* and *isotropic* we can show that

$$\gamma_z(\mathbf{h}) = \sigma_z^2 - C_z(\mathbf{h}).$$

- Therefore, if we know the covariance function, we can calculate the (semi)-variogram.

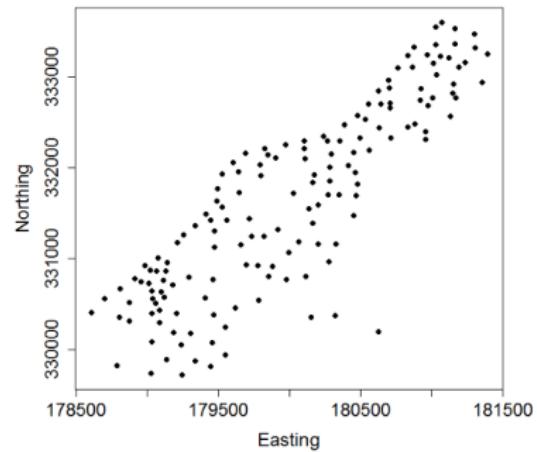
What does the variogram look like?

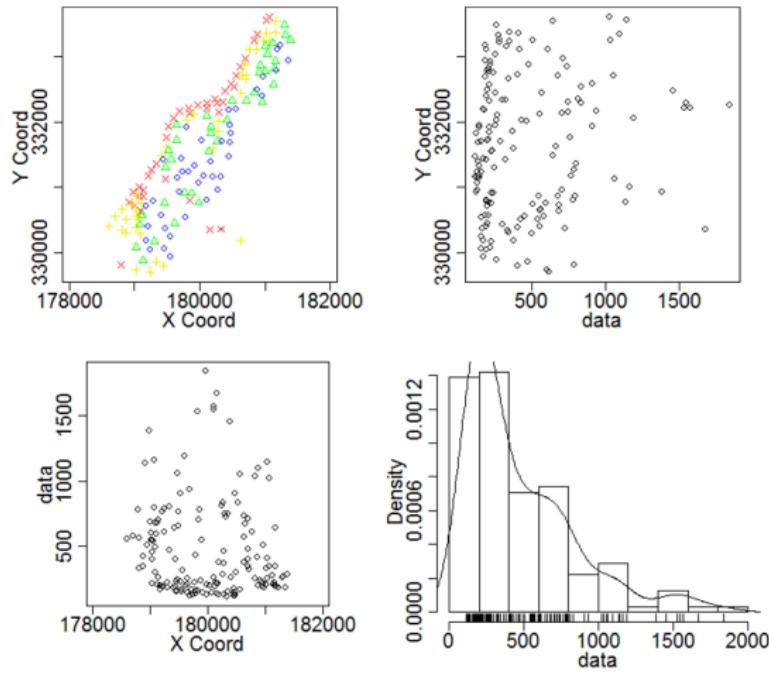
- The plot below shows a typical variogram, with the variance increasing as the lag increases.
- The **sill** is the maximum variance as $h \rightarrow \infty$.
- The **nugget** is the minimum variance as $h \rightarrow 0$.
- The **range** is the distance to the sill.
- Points further apart than the range are assumed to be uncorrelated.



- The variogram is a function of the underlying geostatistical process Z .
- In practice, we only have access to m realisations of this process, and therefore we have to estimate the variogram.
- This is known as the *empirical variogram*.
- We obtain this by computing the semi-variance for all possible pairs of observations: $\gamma(s, t) = 0.5(Z(s) - Z(t))^2$.
- We can then plot the empirical variogram by plotting these semi-variances against their corresponding lags.

- We are interested in soil zinc levels in the flood plains of the River Meuse (in France, Belgium and the Netherlands).
 - Our goal is to model the spatial correlation in the data so that we can predict zinc levels at unsampled locations.
-
- A set of sampling locations were selected.
 - Our data consist of location coordinates and zinc levels.

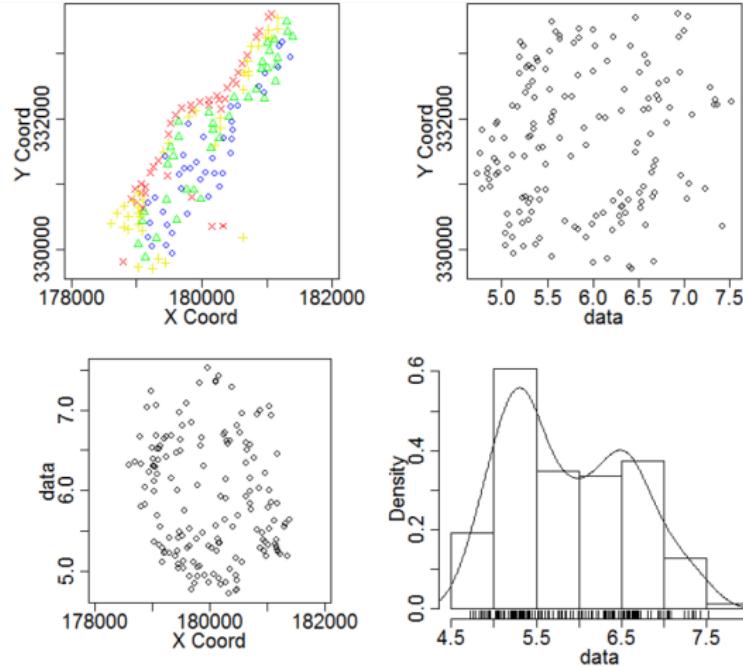




Plots: (a) map colour coded by zinc level, (b) and (c) the zinc levels plotted against the x and y coordinates and (d) a density plot of the zinc levels.



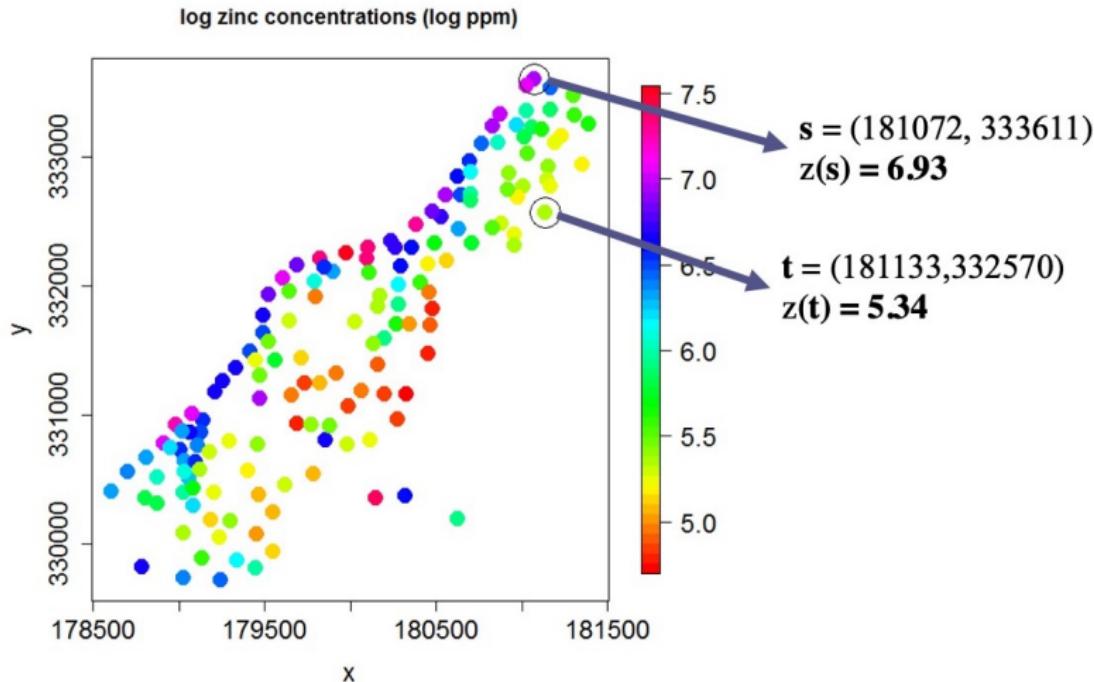
<https://www.menti.com/alws53ba6bzm>



Plots: (a) map colour coded by *log* zinc level, (b) and (c) the *log* zinc levels plotted against the x and y coordinates and (d) a density plot of the *log* zinc levels.

Computing the empirical variogram

- To illustrate how an empirical variogram is computed, consider the two highlighted locations below.



- We can first compute the distance between the two locations using the standard Euclidean distance formula as

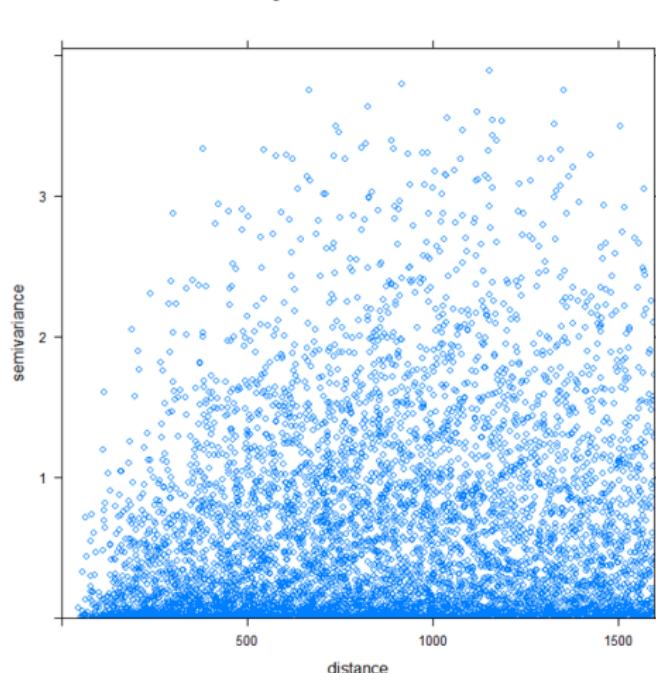
$$h = \sqrt{(181072 - 181133)^2 + (333611 - 332570)^2} = 1042.8$$

- Next, we compute the semi-variance between the points using their observed values as

$$\begin{aligned}\gamma(s, t) &= 0.5(Z(s) - Z(t))^2 \\ &= 0.5(6.93 - 5.34) = 1.27\end{aligned}$$

- We repeat this process for every possible pair of points, and plot h against $\gamma(s, t)$ for each.

- This plot shows the semi-variances for each pair of points.
- Each pair of points has a different distance, making it difficult to use this for prediction.

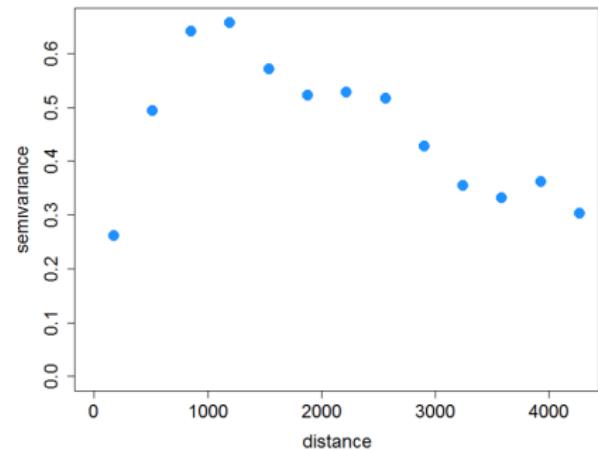
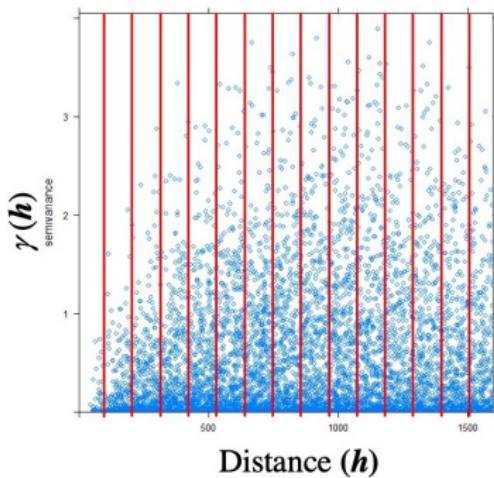


- To make the variogram easier to use and interpret, we divide the distances into a set of discrete bins, and compute the average semi-variance in each.
- We compute this binned empirical variogram as

$$\gamma(\mathbf{h}) = \frac{1}{2N(h_k)} \sum_{(\mathbf{s}, \mathbf{t}) \in N(h_k)} [z(\mathbf{s}) - z(\mathbf{t})]^2$$

- Here, k is the number of bins and $N(h_k)$ is the number of points in the bin with average distance h .
- We then construct a plot of our empirical variogram and use this to estimate the covariance structure.

- The bins are illustrated on the left, and the empirical variogram obtained from them is shown on the right.
- We can start to think about estimating a nugget, sill and range based on this empirical variogram.

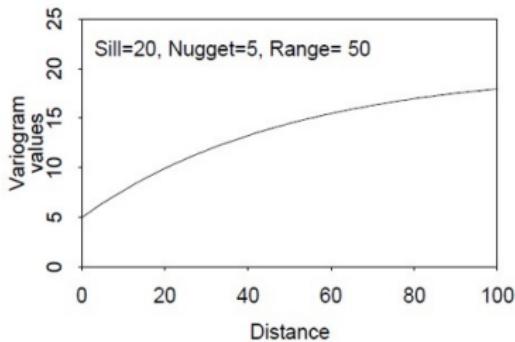


- Once we have computed an empirical variogram, we have to think about fitting a model to it.
- We only have estimates at a small number of lags, we would like a continuous model which explains the changes in dependence structure as h increases.
- Our variogram model needs to have the following properties:
 - Monotonically increasing
 - A constant maximum (sill)
 - A positive intercept (nugget)
- Several models exist which satisfy these criteria.

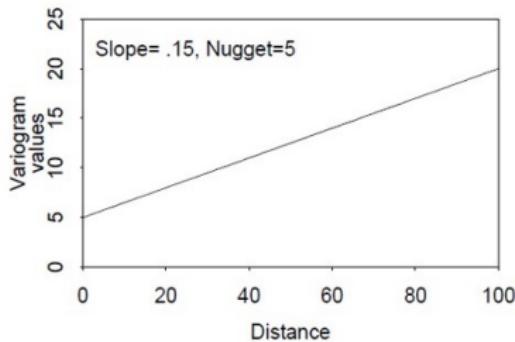
- According to Webster and Oliver (2007), choosing variogram models is one of the most controversial topics in geostatistics.
- We often choose the right model based on mathematical criteria such as least squares, maximum likelihood or AIC.
- One of the more popular approaches, proposed by Cressie (1985), is to use weighted least squares, where the weights are based on the number of observations within each ‘bin’.
- The outcome of this is that more weight is given to the lags which have been estimated with more data points, which are usually the shorter lags.

Examples of variogram models

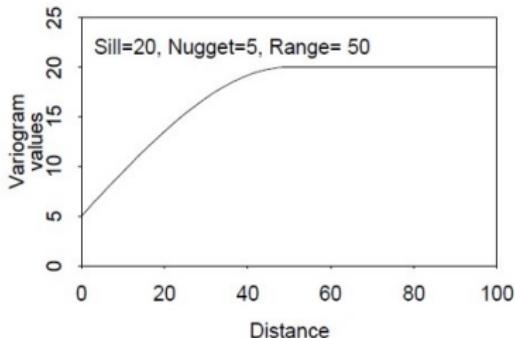
Exponential Variogram



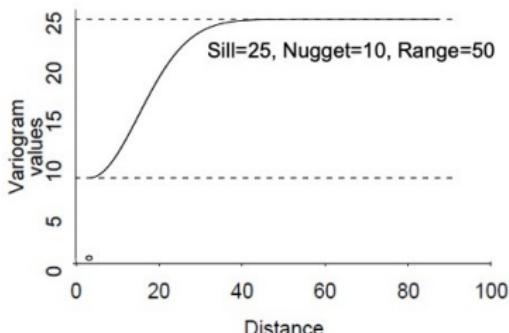
Linear Variogram



Spherical Variogram



Gaussian Variogram



Random $\gamma(h) = \begin{cases} 0, & \text{if } h = 0 \\ c, & \text{otherwise} \end{cases}$

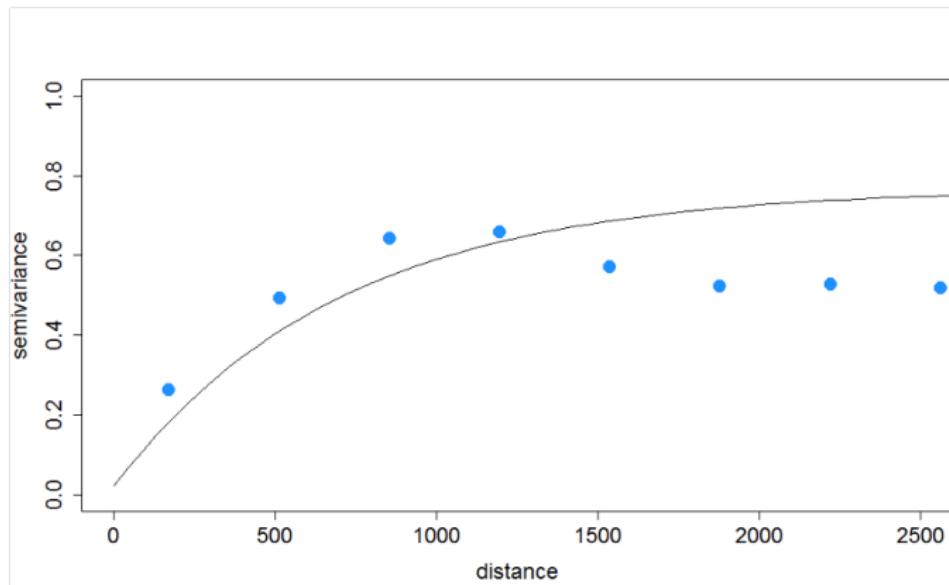
Spherical $\gamma(h) = \begin{cases} c \cdot 1.5 \frac{h}{a} - 0.5 \left(\frac{h}{a}\right)^3, & \text{if } h \leq a \\ c, & \text{if } h \geq a \end{cases}$ *a is the range of influence.*

Exponential $\gamma(h) = \begin{cases} \tau^2 + \sigma^2(1 - \exp(-h/\phi)) & \text{if } h > 0 \\ 0, & \text{if } h = 0 \end{cases}$

Gaussian $\gamma(h) = \begin{cases} \tau^2 + \sigma^2(1 - \exp\{(-h/\phi)^2\}) & \text{if } h > 0 \\ 0, & \text{if } h = 0 \end{cases}$

- The key features of our variogram are represented by the following model parameters
 - $\tau^2 > 0$ is the **nugget**.
 - $\sigma^2 > 0$ is the **partial sill**.
 - $\phi > 0$ is the **range parameter**.
- Note that ϕ is not the range itself, but rather a parameter which controls how quickly the covariance decays towards zero (or the variogram increases to the sill).
- A smaller value of ϕ means the covariance function decays to zero quickly, a larger value means it decays to zero more slowly.

- We can now fit a variogram model to the empirical variogram obtained from the Meuse river example.
- The black solid line below shows an exponential variogram model (fitted using the `geoR` package).



- Once we have estimated a variogram to account for our spatial autocorrelation, we can start to think about making predictions.
- Spatial prediction is the process of predicting the value of our variable of interest at an unobserved location s_0 .
- As with any statistical prediction, we use what we know about our observed data, including their values, how far our unobserved location is from them, and our variogram.
- There are many methods for spatial prediction, including regression modelling, distance weighted interpolation and an approach known as **kriging**.

- Kriging is an approach named after its inventor D. G. Krige, who worked in the mining industry in South African in the 1950s.
- He used this approach to understand the spatial pattern of mineral resources.
- It is a relatively simple and theoretically appealing approach, and is therefore incredibly popular for geostatistical prediction.
- Kriging interpolates between previously observed locations in order to predict at new locations.

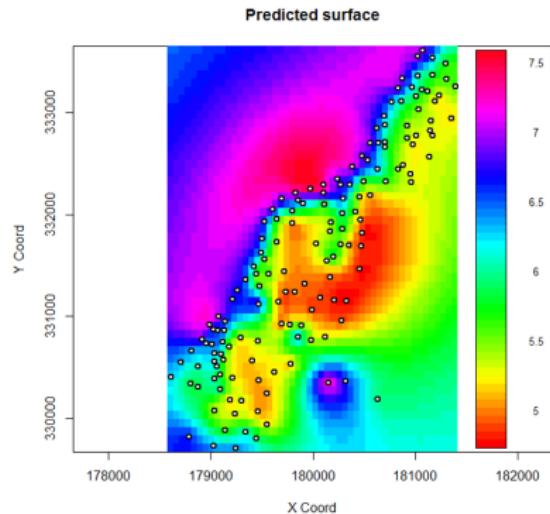
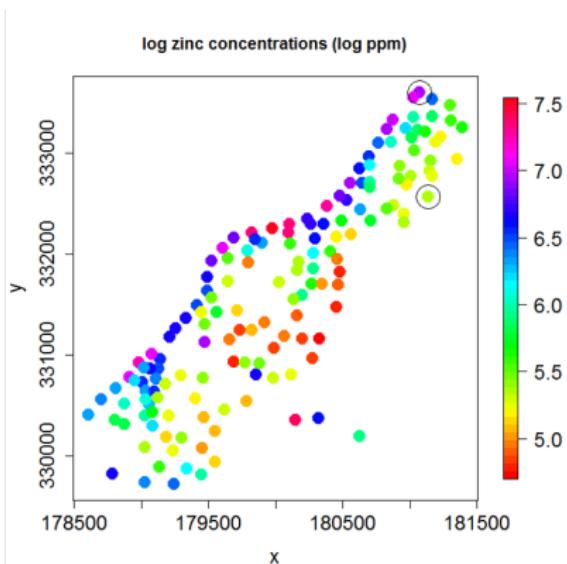
- Ordinary kriging is a form of kriging which assumes that the overall mean and variance of the region is constant.
- Predictions at unsampled locations are made using a weighted average of the observations.

$$z^*(\mathbf{s}_0) = \sum_{i=1}^m \lambda_i z(\mathbf{s}_i)$$

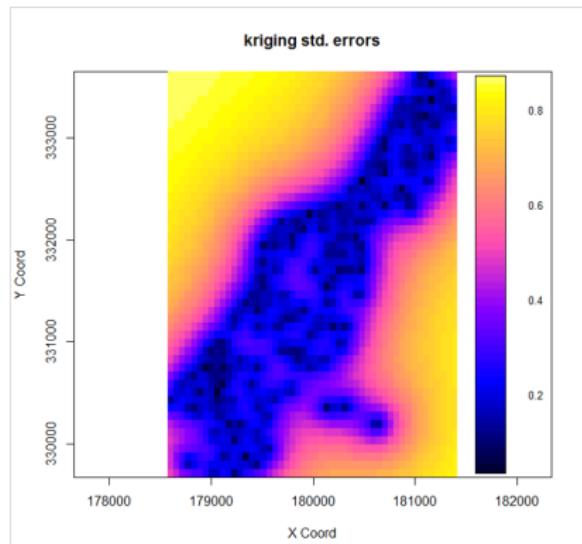
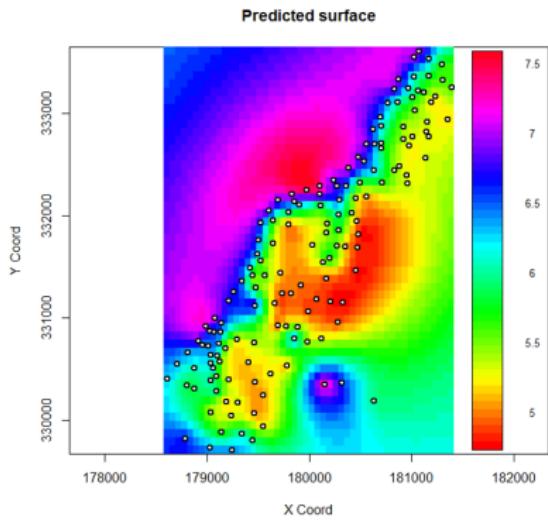
- The weights λ_i can be estimated in a number of different ways, and are commonly based on the variogram.
- The weights are therefore usually proportional to the distance between the observed and new locations.

- There are other types of kriging such as *universal kriging* and *block kriging* which are more appropriate for different structures of spatial data.
- Whichever method is used, we can obtain a set of predictions over a grid to allow us to plot a map of the predicted spatial surface for our variable of interest.
- Our predictions will be of better quality in areas where there are lots of observed values, and poorer quality when we have less nearby data available.
- Therefore we can also generate a map of the uncertainty across the region.

- The surface map on the right is generated by using kriging to make predictions over a fine grid.
- Notice the ‘edge effects’ - a gradual smoothing towards the mean as you move away from any actual data in the top left and bottom right.



- The map on the left shows the uncertainties associated with our estimated surface.
- There is lower uncertainty in the areas with lots of observed data, and higher uncertainty as we move away from the observed data.



- When working with geostatistical data, we first have to understand what sort of spatial trend, if any, is present in our data.
- Next, we plot the empirical (semi-)variogram and choose a suitable model structure.
- We then identify the parameters (sill, nugget, range) of this chosen model structure.
- Finally, we use this model to obtain weights which can be used for interpolation via, for example, kriging, and use this to estimate the surface on a map.



<https://www.menti.com/alws53ba6bzm>

Temporary page!

\LaTeX was unable to guess the total number of pages correctly. As there was some unprocessed data that should have been added to the final page this extra page has been added to receive it.

If you rerun the document (without altering it) this surplus page will go away, because \LaTeX now knows how many pages to expect for this document.