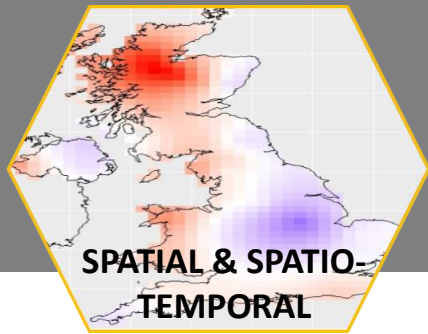


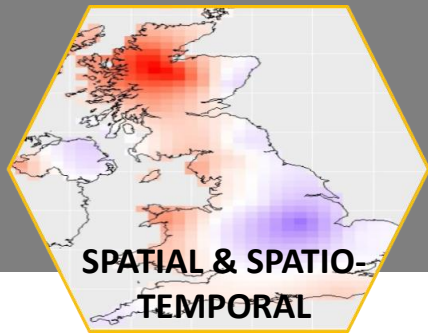
SPATIAL & SPATIO-TEMPORAL ANALYSIS – PART I



SPATIAL DATA ANALYSIS

CONTENT & LEARNING OUTCOMES

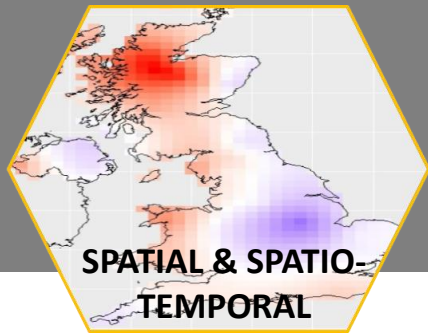
- Spatial data types
- Introduction to Geostatistics (focus on Gaussian data)
- Why 'go Bayesian' for spatial/spatio-temporal problems
- Examples
- In practice (software/packages)



SPATIAL DATA ANALYSIS

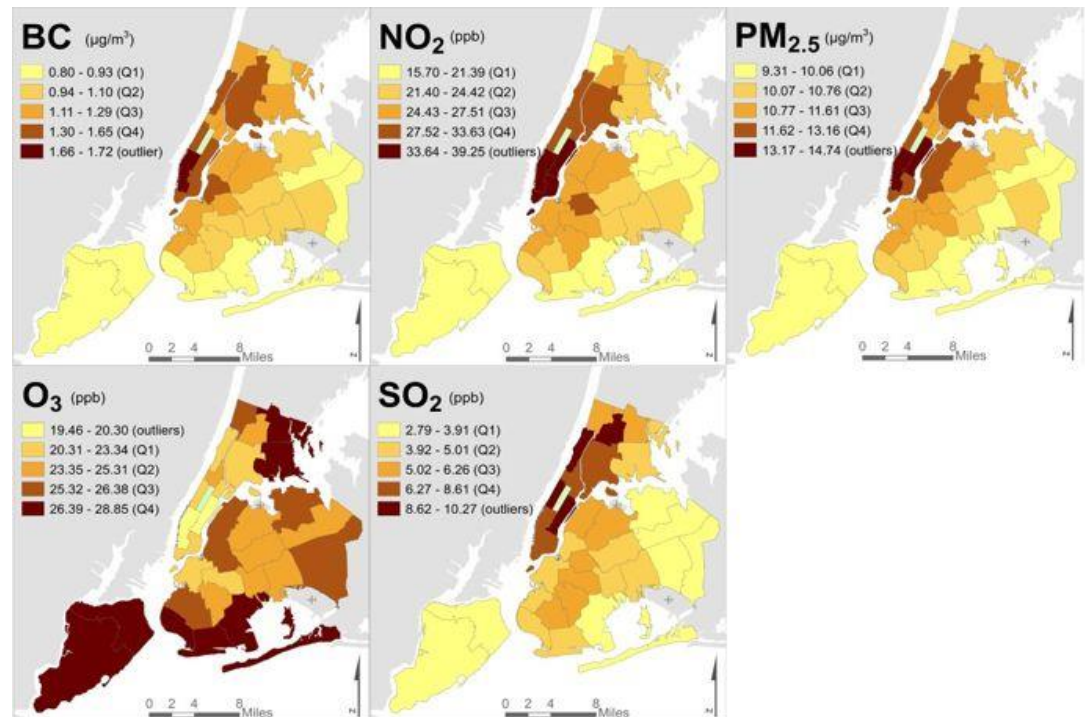
SPATIAL DATA '101'

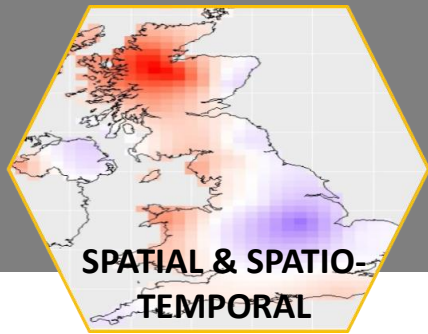
- Data collected in space
- i.e. Sampling location is known and relevant
 - Space is of primary interest – we want to analyse and describe the spatial structure, how the data vary over space
 - Spatial structure is nuisance – we need to account for it in order to draw proper inference from our models (spatial autocorrelation)
 - Often data collected from more similar locations are more similar so are not independent
 - Upscaling



SPATIAL DATA FORMS

- Areal data – partitioned region with discrete spatial units, each with one value

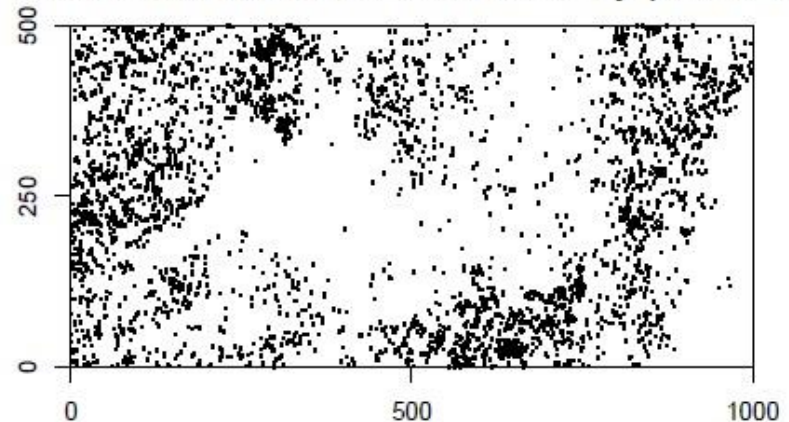


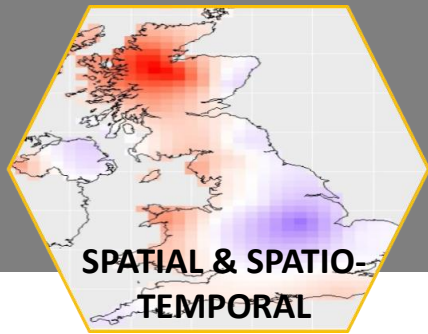


SPATIAL DATA FORMS

- Areal data
- Spatial point pattern data – Pattern formed by location of objects/events. Is the spatial pattern of points random or clustered/structured? Are the patterns determined by covariates?

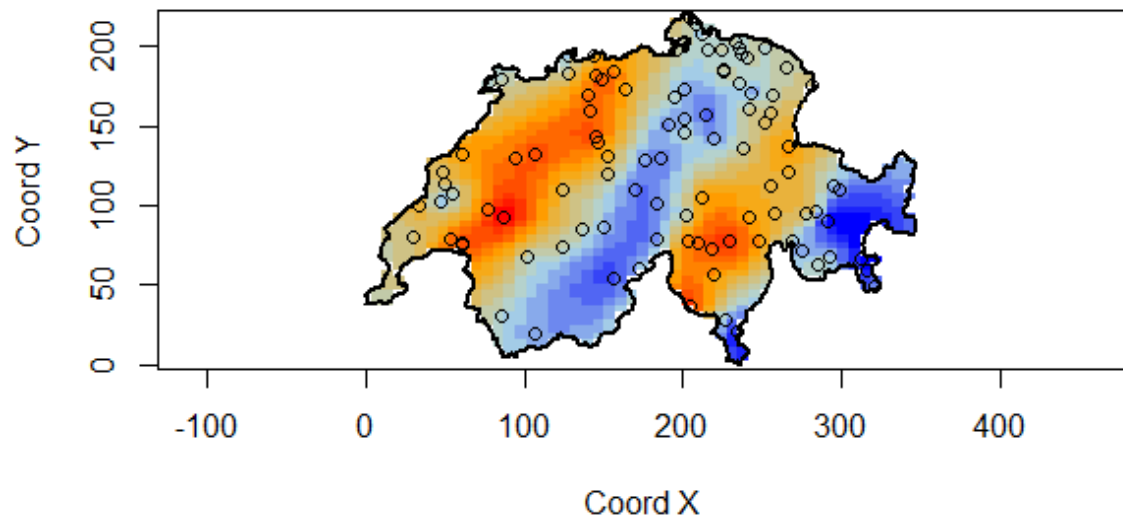
- locations of rainforest trees on study plot in Panama

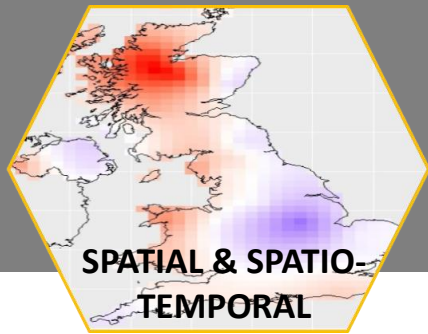




SPATIAL DATA FORMS

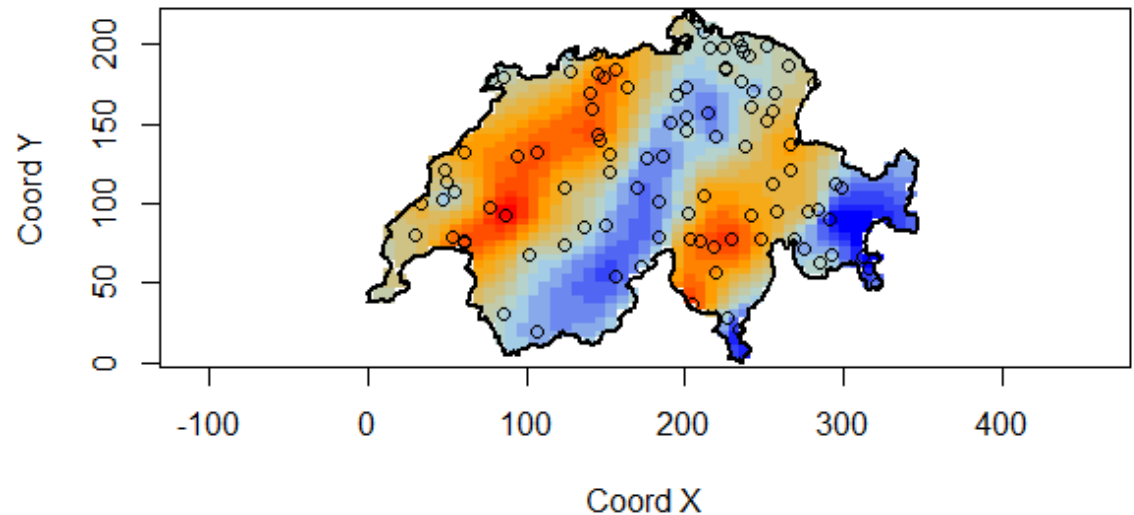
- Areal data
- Spatial point pattern data
- Geostatistical data – spatially continuous phenomena, based on observations at a finite number of locations



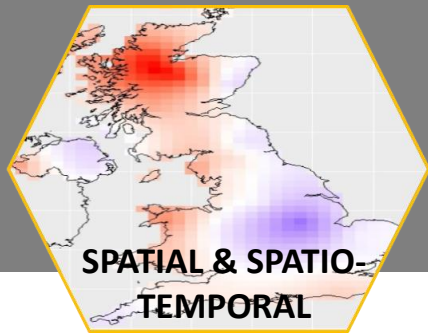


SPATIAL DATA FORMS

- Areal data
- Spatial point pattern data
- Geostatistical data

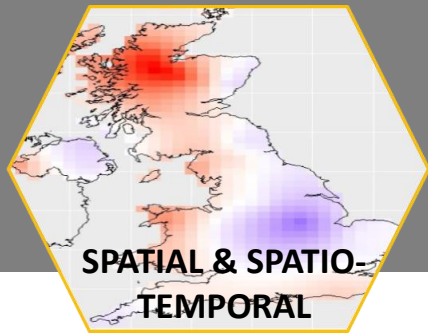


- Sometimes this is not our end-point -> covariates for other environmental/ecological processes



GEOSTATISTICS

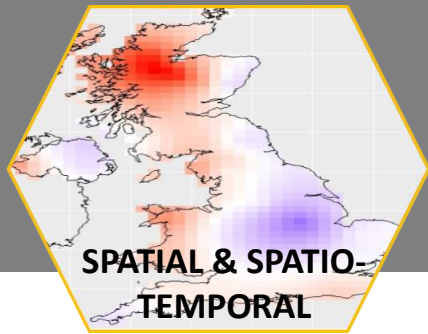
- SAMPLES drawn from a CONTINUUM in space
- Often a goal is to PREDICT values in intervening spaces
- Without BIAS and MINIMISING UNCERTAINTY



GEOSTATISTICS

SPATIAL DEPENDENCE or CORRELATION

- **Correlations** between sites are a function of **distance**
- Typically, geostatistical data will display **positive** correlation
- The closer two observations are the more similar their values are likely to be
- Arises because variables of interest being affected by other **unmeasured processes** which are themselves spatially correlated
- For example:
 - Air pollution
 - Soil nutrients
 - Your examples.....

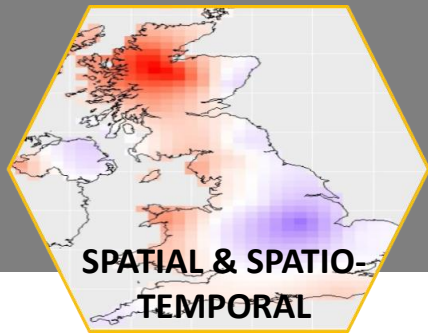


WHY 'GO BAYESIAN'...

...FOR GEOSTATISTICAL ANALYSIS

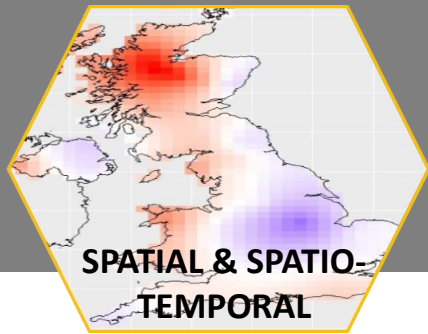
- Correctly allows for variation in the parameters -
Parameters of correlation function – ***random***, not fixed
(have an associated probability distribution)
- Uncertainty intervals are easy to obtain for all parameters,
not just regression parameters
- Appropriate propagation of uncertainty means prediction
intervals will be wider

...WE WILL REVISIT!



A NOTE ON SOFTWARE

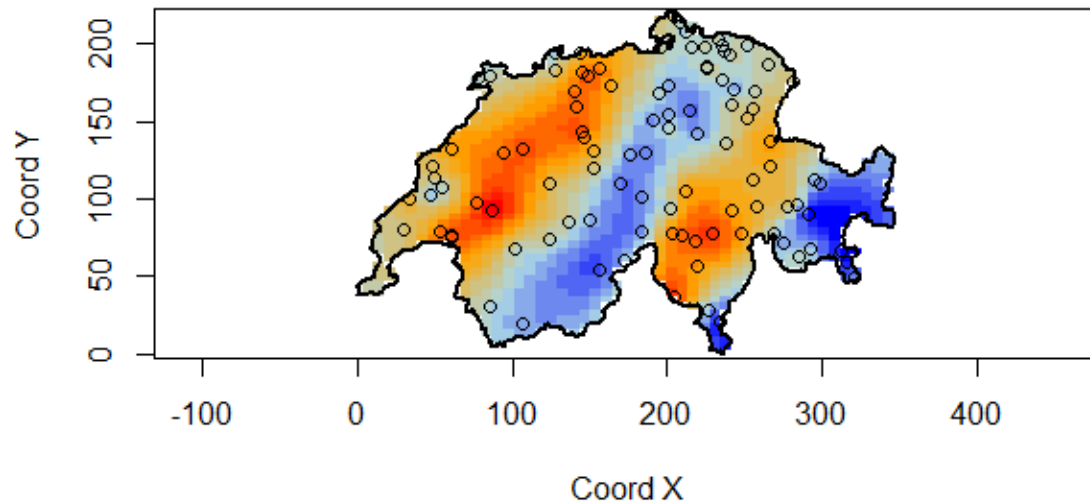
- **geoR**, and **geoRglm** for non-Gaussian data
- **SpBayes**
- **R-INLA/inlabru**
- **OpenBUGS**
- **JAGS**



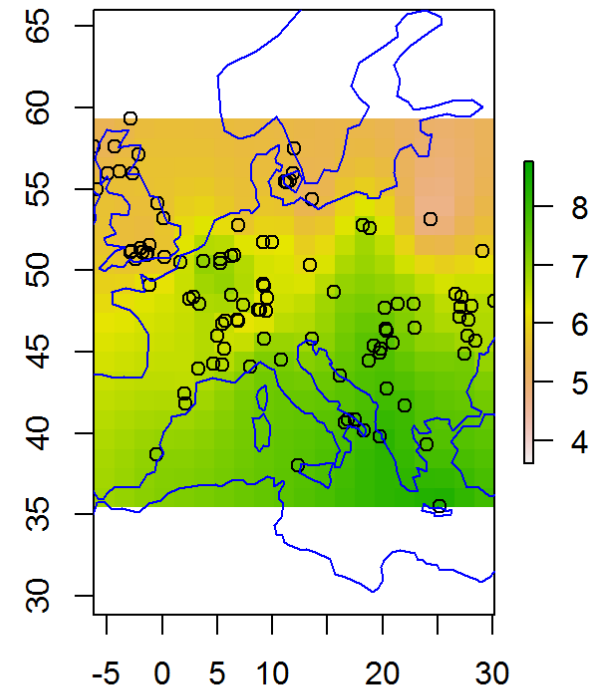
TODAY'S PRACTICALS

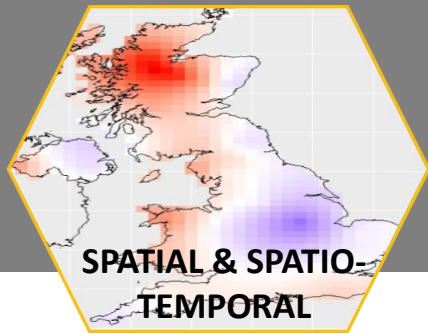
In geoR....

Swiss Rainfall Data



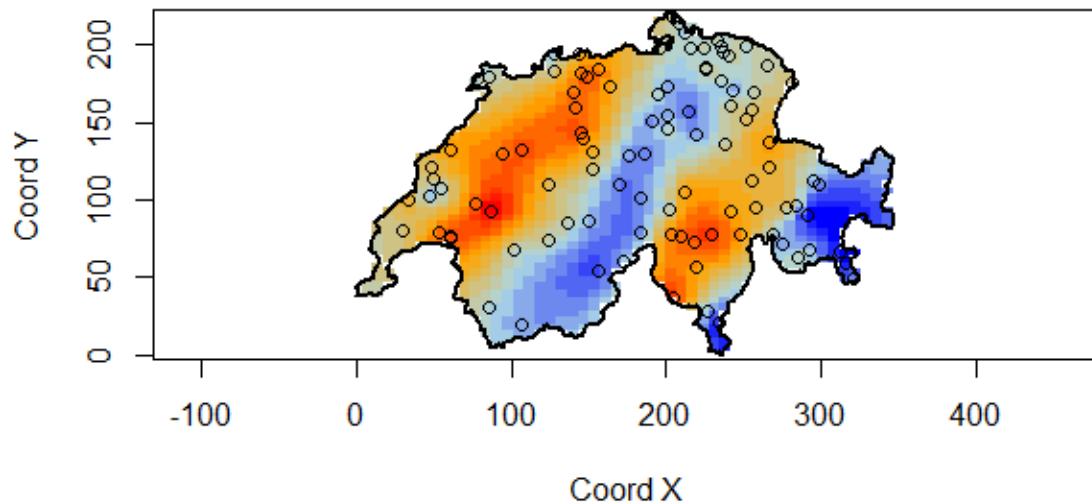
Spread of Agriculture





GEOSTATISTICS

Swiss Rainfall Data



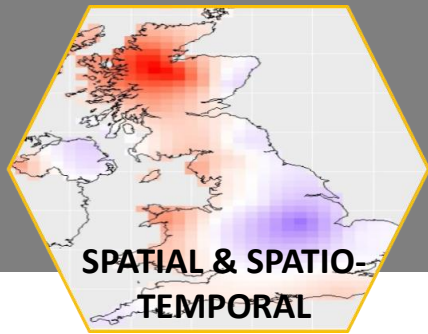
Stochastic (random) process
- Variable of interest, Z , at
locations set (s) within total
space D

$$\{Z(\mathbf{s}) : \mathbf{s} \in D\}$$

Locations s at which data
could occur vary
continuously over D .
But data are observed at a
finite number of locations,
denoted by:

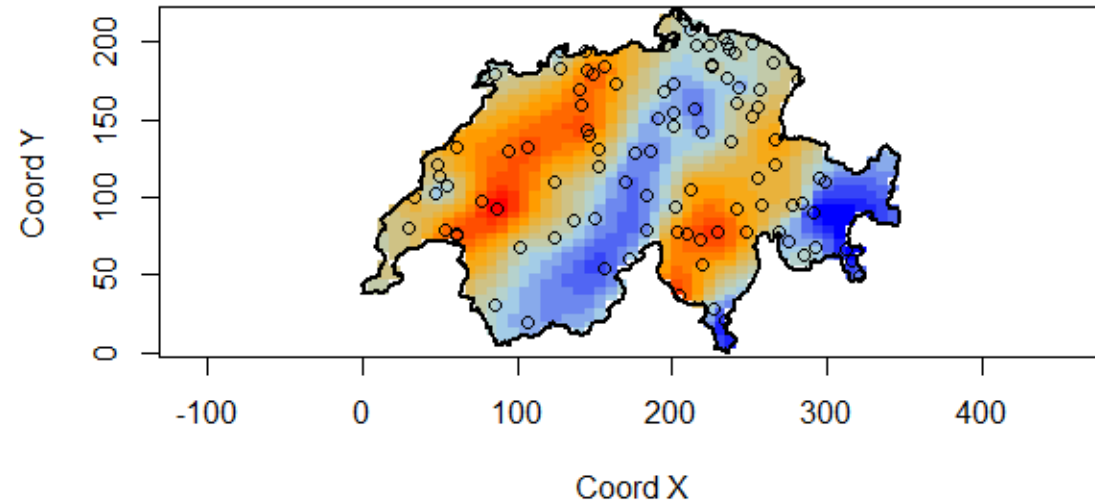
$$\mathbf{z} = \{z(\mathbf{s}_1), \dots, z(\mathbf{s}_m)\}$$

i.e. a particular realisation of
random variables $Z(s)$



GEOSTATISTICS

Swiss Rainfall Data



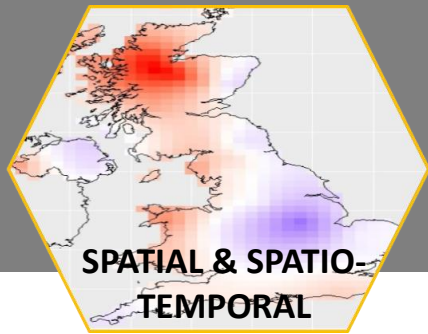
Data are observed at m locations:
What drives variation in Z ?
Often want to predict the unknown stochastic process at locations where we have not sampled – i.e. produce a map across domain D .

Stochastic (random) process:
Variable of interest, Z , at locations (s) within total space D

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

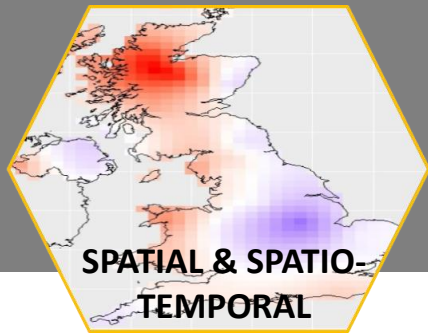
Locations s at which data *could* occur vary continuously over D .
But data are observed at a finite number of locations, denoted by:

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



STATIONARITY

- Data are observed at m locations but we often want to predict the unknown stochastic process at locations where we have not sampled – i.e. produce a map across domain D .
- BUT we only have one REALISATION – inference requires many realisations
- We must make some assumptions to simplify
- STATIONARITY – such that each *observation* can be treated as a random variable
 - Strict stationarity – all characteristics of the random function remain the same (not practicable)
 - Weak/Intrinsic stationarity – certain moments are invariant, whereas others are allowed to vary
- Put simply, the process $Z(s)$ has the same degree of variation from place to place
- These assumptions underlie the theoretical variogram

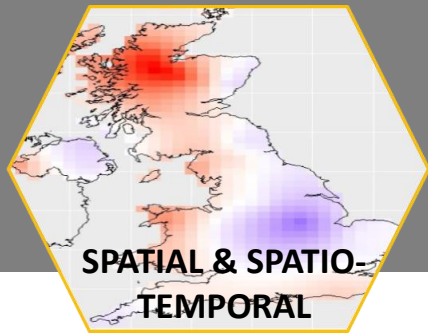


RANDOM PROCESS MODEL

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

Mean of the process μ , and random quantity $\varepsilon(\mathbf{s})$ with mean 0 and covariance $C(h)$, where h is the separation between samples in distance and direction.

$$C(\mathbf{h}) = E[Z(\mathbf{s})Z(\mathbf{s} + \mathbf{h}) - \mu^2]$$



THE VARIOGRAM

Spatial association as a function of separation distance

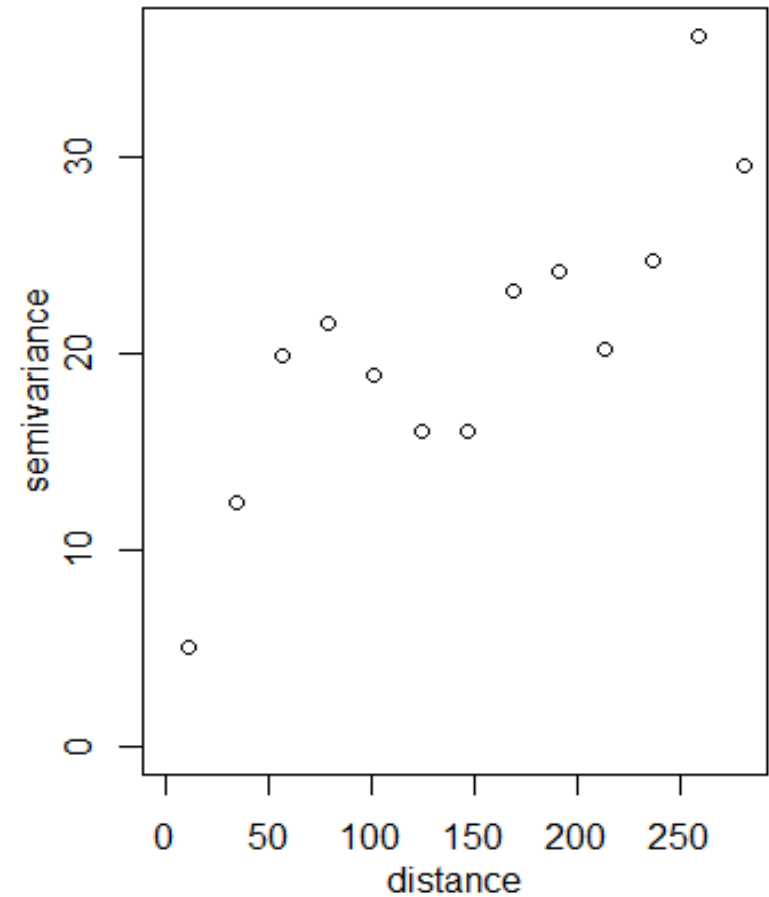
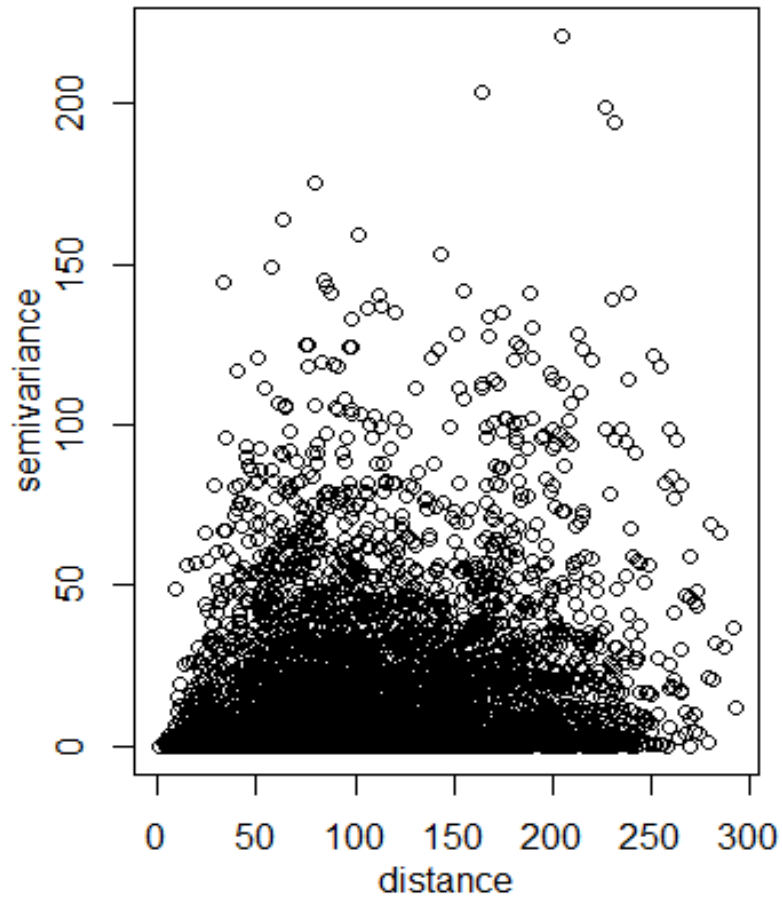
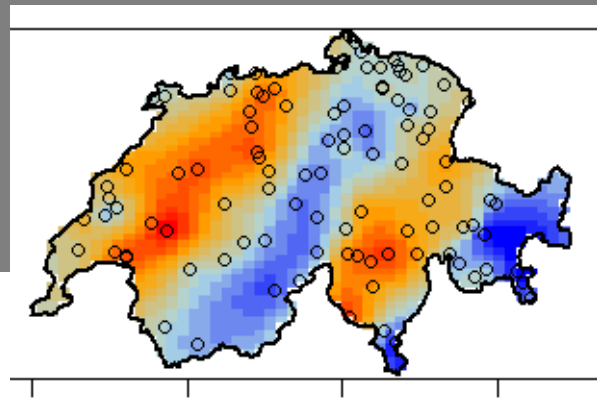
The semi-variance is a function denoted by:

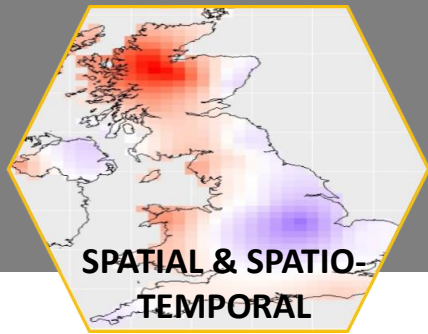
$$\gamma_z(\mathbf{h}) = \frac{1}{2} \text{Var}[Z(\mathbf{s}) - Z(\mathbf{s} + \mathbf{h})]$$

Semi-variance depends only on h (separation) – as a function of h , it is *the variogram*

Variance in difference in the process Z between two locations – when this is small, locations are spatially correlated

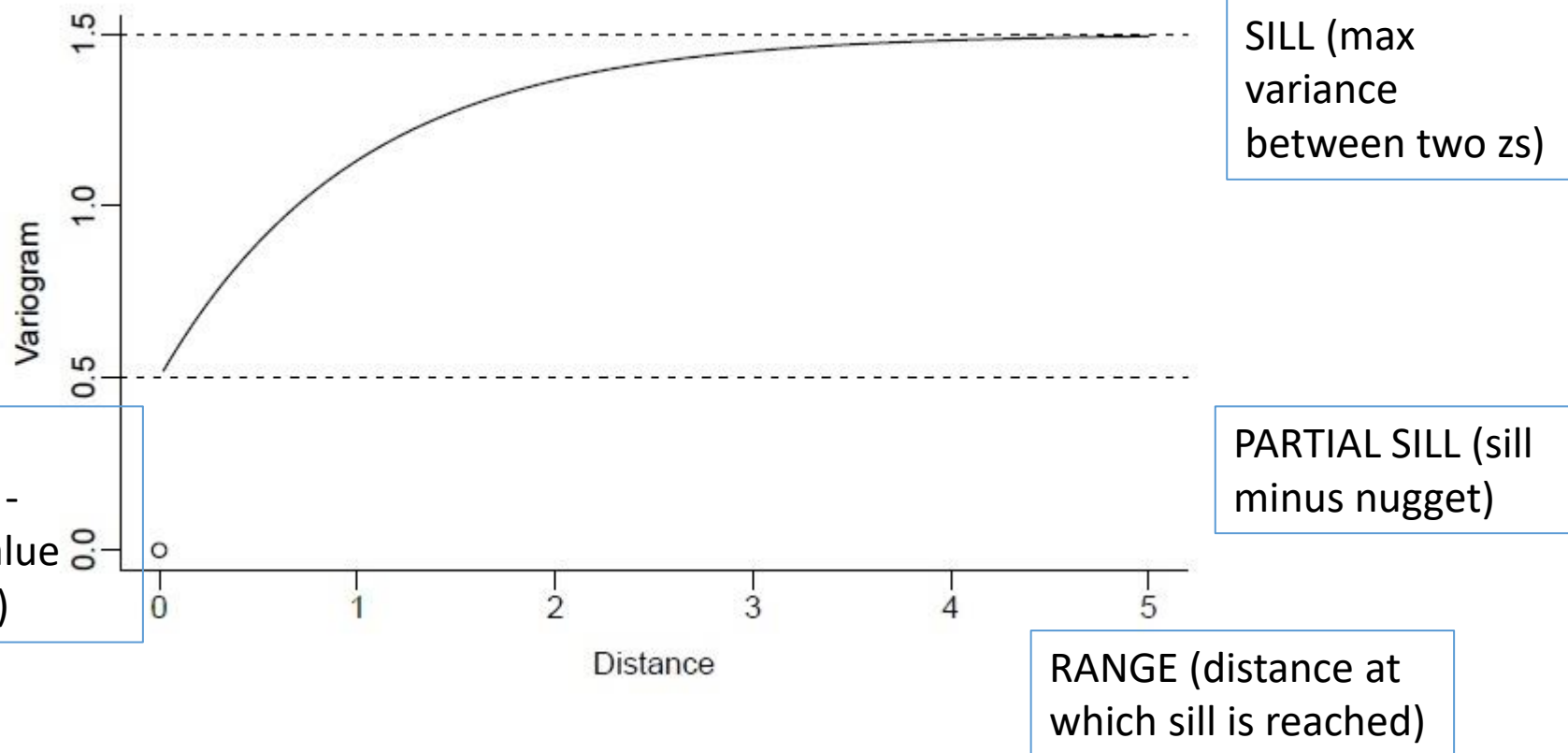
EXPERIMENTAL VARIOGRAM

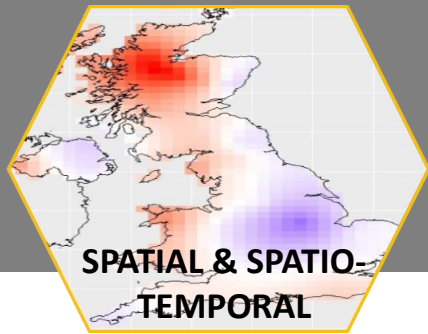




THE VARIOGRAM

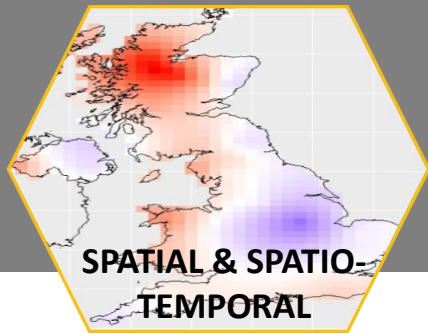
Theoretical shape of the variogram





PARAMETRIC MODELS OF THE VARIOGRAM

- Features:
 - Nugget – τ^2 [τ^2] > 0
 - Partial sill – σ^2 [σ^2] > 0
 - Range parameter (rate at which covariance decays to zero) – ϕ [ϕ] > 0



PARAMETRIC MODELS OF THE VARIOGRAM

- Some commonly used correlation functions

Definition 5.3.2 — Exponential model.

$$\gamma(h) = \begin{cases} 0 & h = 0 \\ c_0 + c_s [1 - \exp(-\frac{h}{a})] & h > 0 \end{cases} \quad (5.16)$$

Floch 2018 *Handbook of Spatial Statistics*

Definition 5.3.3 — Gaussian model.

$$\gamma(h) = \begin{cases} 0 & h = 0 \\ c_0 + c_s [1 - \exp(-(\frac{h}{a})^2)] & h > 0 \end{cases} \quad (5.17)$$

Definition 5.3.4 — Power model.

$$\gamma(h) = \begin{cases} 0 & h = 0 \\ c_0 + bh^p & h > 0 \end{cases} \quad (5.18)$$

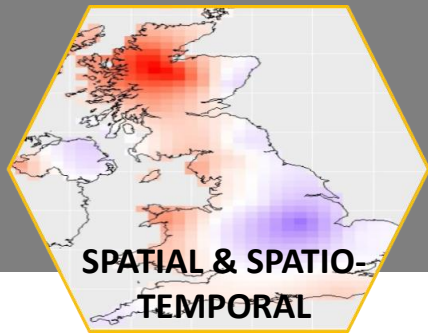
Definition 5.3.5 — Matern model.

$$\gamma(h) = \begin{cases} 0 & h = 0 \\ c_s \left[1 - \frac{\frac{h}{a}}{2^{\alpha-1}\Gamma(\alpha)} K_{\alpha} \left(\frac{h}{a} \right) \right] & h > 0 \end{cases} \quad (5.19)$$

where Γ refers to the gamma function and K_{α} , the modified Bessel of the second kind of parameter α .

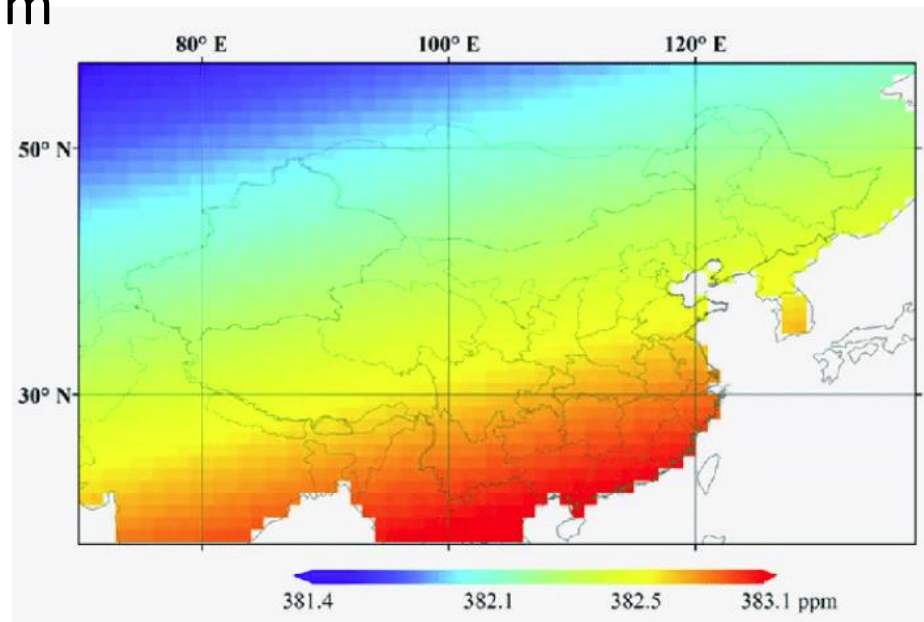
- Caution with Gaussian – can lead to bizarre predictions (Wackernagel 2003)
- Choice of function -> influence results**





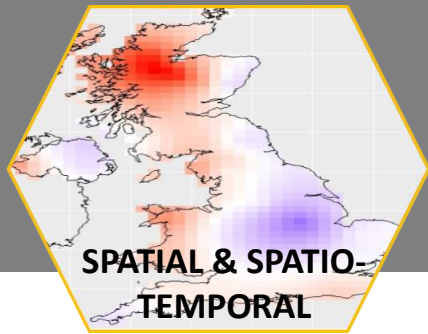
TRENDS IN SPATIAL DATA

- Recall stationarity assumption....
- Gradual variation/large-scale variation
- Modelled using covariates or coordinates
- Then, in essence, just modelling the *residual* spatial variability with the variogram

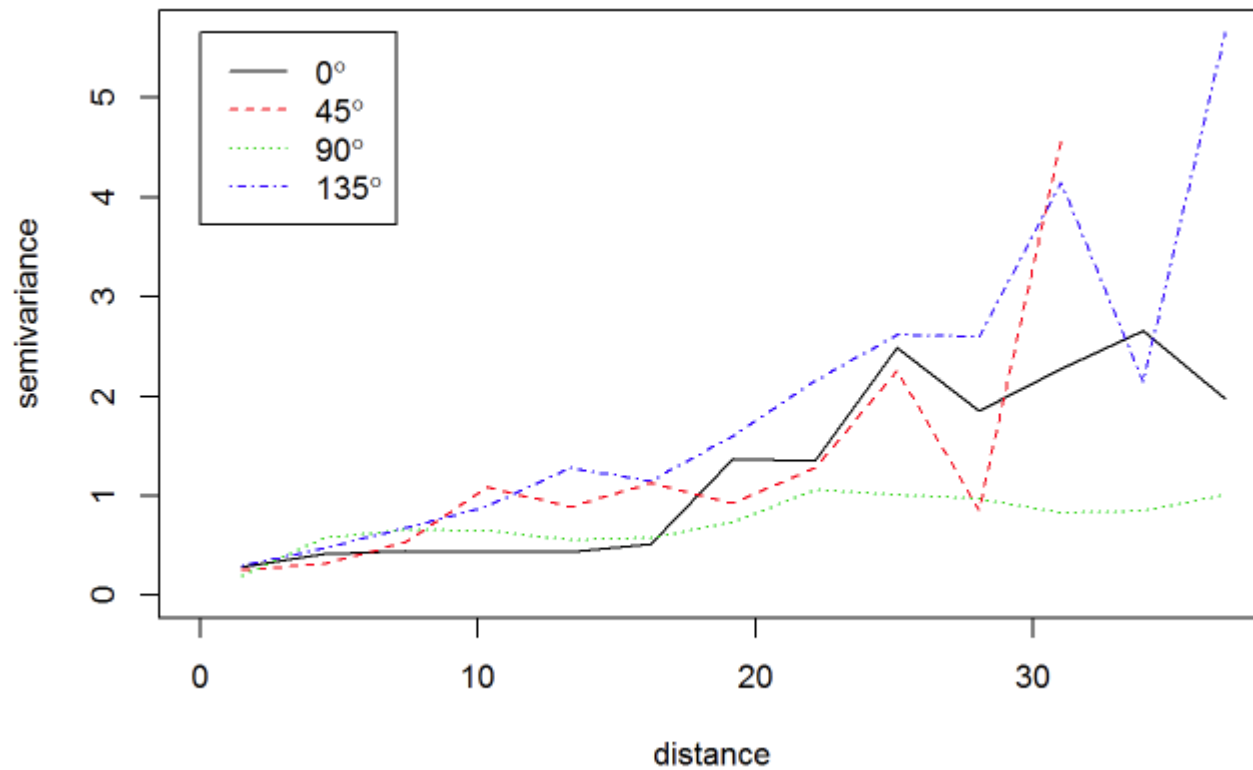


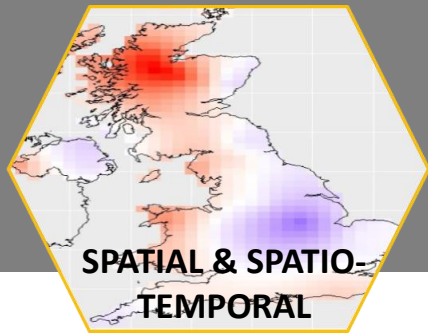
Zeng et al. 2014 *Transactions on Geoscience and Remote Sensing*

Linear spatial trend surface of Xco₂ in the study area derived from ACOS-GOSAT data in one month of September 2009.



ISOTROPY





BAYESIAN GEOSTATISTICAL MODEL

Data:

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

Gaussian model with likelihood:

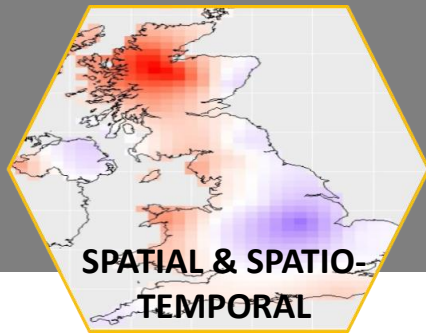
$$\mathbf{z} \sim N(X\boldsymbol{\beta}, \Sigma(\boldsymbol{\theta}))$$

Where the mean function $X\boldsymbol{\beta}$ is a linear combination of known covariates (X) and $\Sigma(\boldsymbol{\theta})$ is $m \times m$ covariance matrix for the m observations, determined by a stationary and isotropic covariance function.

The unknown parameters are $\boldsymbol{\Theta} = (\boldsymbol{\beta}, \sigma^2, \varphi, \tau^2, \kappa^2)$

Joint distribution:

$$f(\boldsymbol{\Theta}) = f(\boldsymbol{\beta}, \sigma^2, \varphi, \tau^2, \kappa^2)$$



BAYESIAN GEOSTATISTICAL MODEL

Data:

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

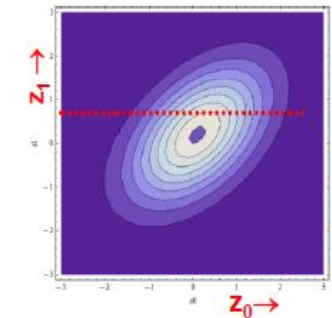
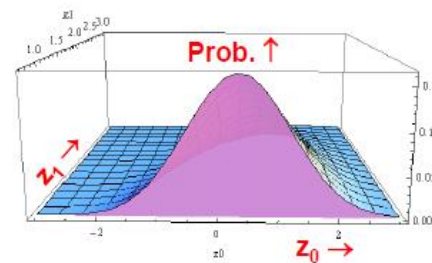
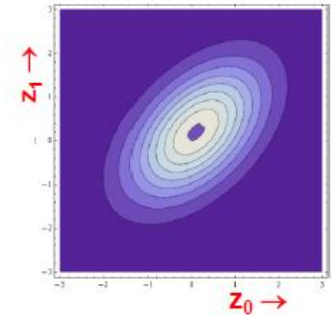
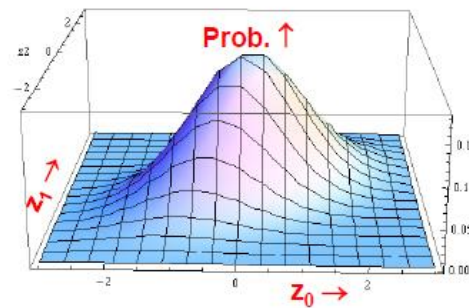
Gaussian model with likelihood:

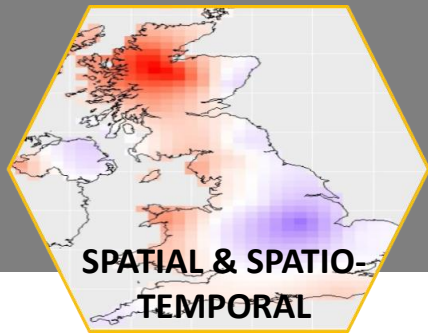
$$\mathbf{z} \sim N(X\boldsymbol{\beta}, \Sigma(\boldsymbol{\theta}))$$

Where the mean function $X\boldsymbol{\beta}$ is a linear (X) and $\Sigma(\boldsymbol{\theta})$ is $m \times m$ covariance matrix determined by a stationary and isotropic
Unknown parameters $\boldsymbol{\Theta} = (\boldsymbol{\beta}, \sigma^2, \varphi, \tau^2, \kappa^2)$,
Joint distribution:

$$f(\boldsymbol{\Theta}) = f(\boldsymbol{\beta}, \sigma^2, \varphi, \tau^2, \kappa^2)$$

Multivariate Gaussian distribution





BAYESIAN GEOSTATISTICAL MODEL

Data:

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

Gaussian model with likelihood:

$$\mathbf{z} \sim N(X\boldsymbol{\beta}, \Sigma(\boldsymbol{\theta}))$$

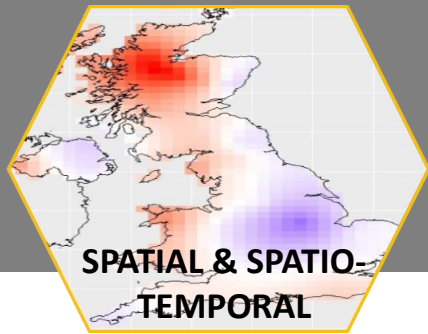
Where the mean function $X\boldsymbol{\beta}$ is a linear combination of known covariates (X) and $\Sigma(\boldsymbol{\theta})$ is $m \times m$ covariance matrix for the m observations, determined by a stationary and isotropic covariance function.

The unknown parameters are $\boldsymbol{\Theta} = (\boldsymbol{\beta}, \sigma^2, \varphi, \tau^2, \kappa^2)$

Joint distribution:

$$f(\boldsymbol{\Theta}) = f(\boldsymbol{\beta}, \sigma^2, \varphi, \tau^2, \kappa^2)$$

DATA + PRIORS \rightarrow POSTERIOR



PLAUSIBLE PRIORS

$$\beta \sim N(\mu_\beta, V_\beta)$$

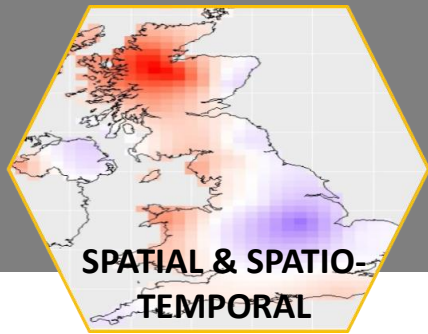
Take on any value

$$\sigma^2, \kappa^2 \sim \text{Uniform}(0, * \text{large}) \text{ or}$$
$$\sigma^2, \kappa^2 \sim \text{Inverse-Gamma}(a, b)$$

Must be positive
Inverse-Gamma conjugate for
variance. Choice of (a,b) under
research/debate.

$$\varphi \sim \text{Uniform}(c, d)$$

Chosen over likely range of
correlations

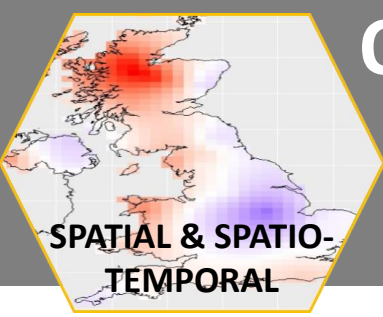


WHY 'GO BAYESIAN'...

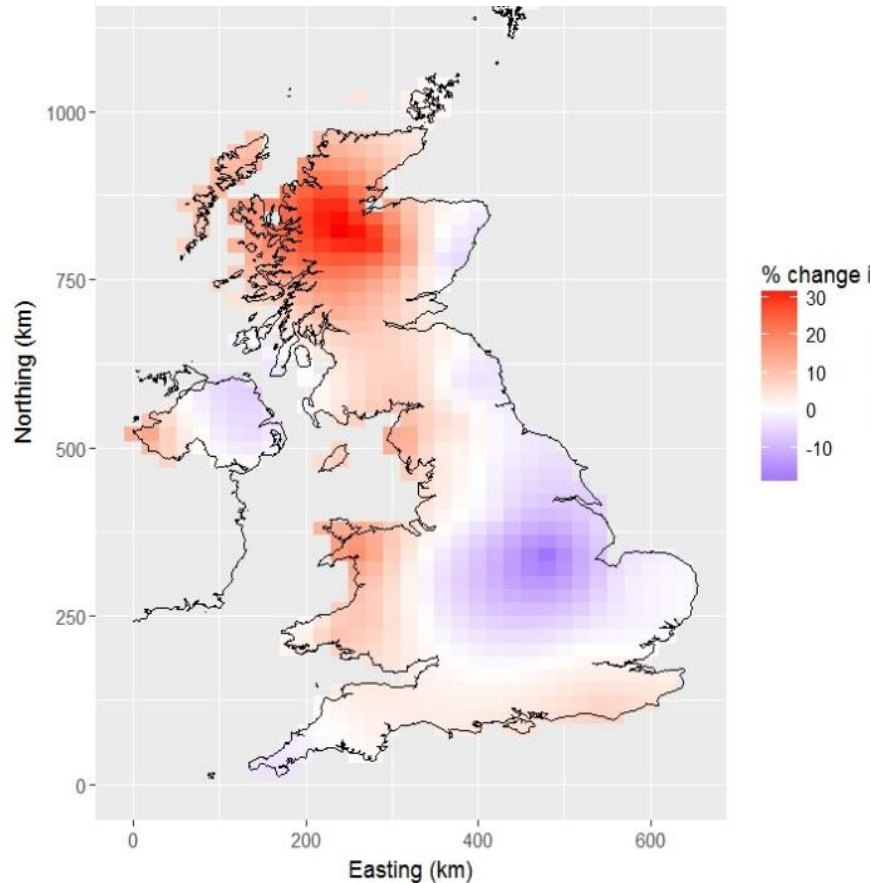
...FOR GEOSTATISTICAL ANALYSIS

- Correctly allows for variation in the parameters -
Parameters of correlation function – *random*, not fixed
(have an associated probability distribution)
- Uncertainty intervals are easy to obtain for all parameters,
not just regression parameters
- Appropriate propagation of uncertainty means prediction
intervals will be wider

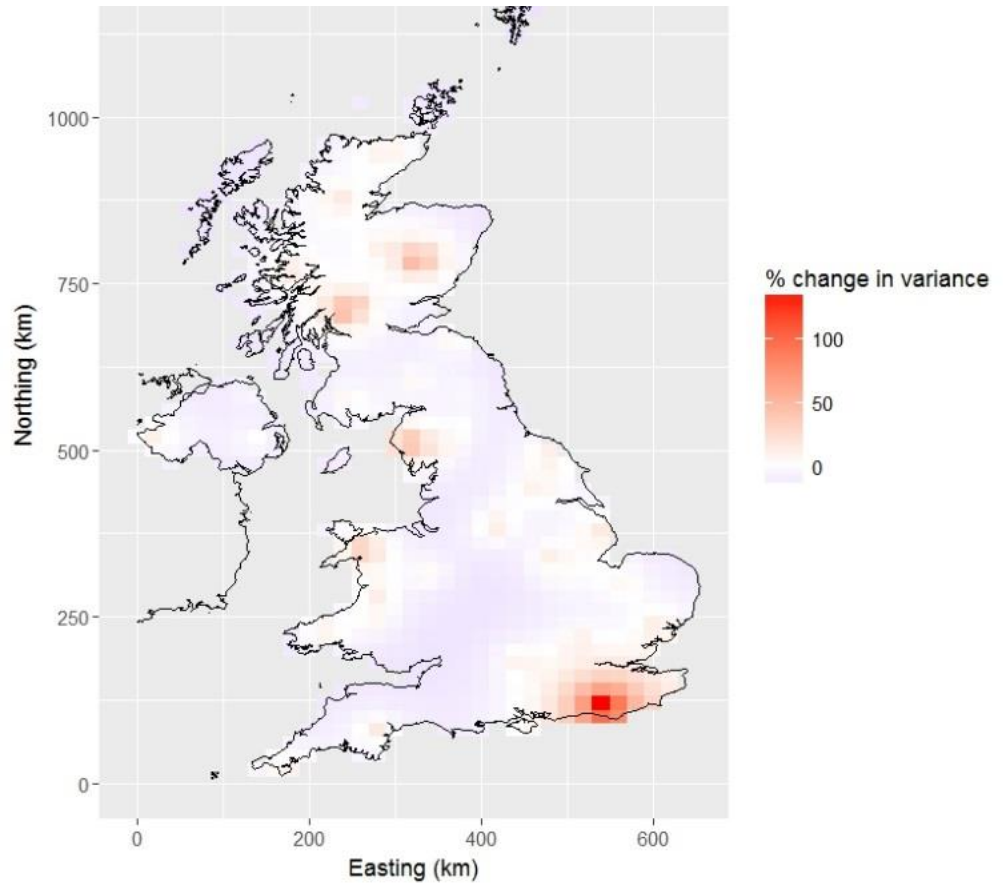
Q: What is the impact of losing 6 sites from an air pollution monitoring network?



Change in estimated ammonium concentration (mean)



Change in variance of ammonium predictions

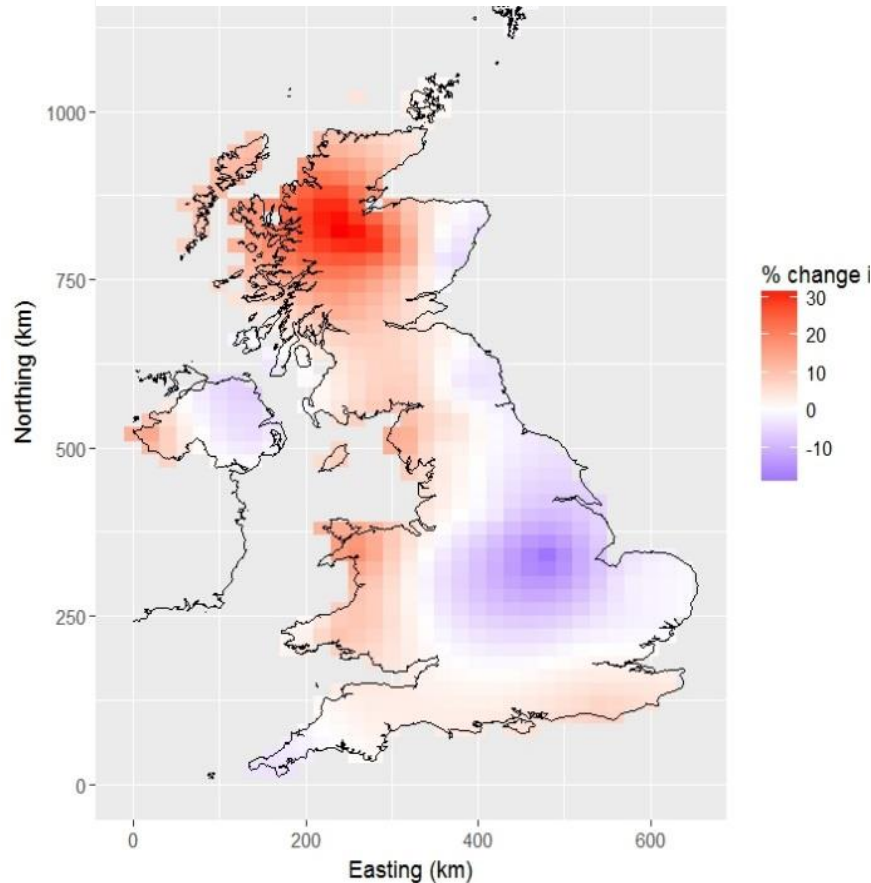


Q: What is the impact of losing 6 sites from an air pollution monitoring network?

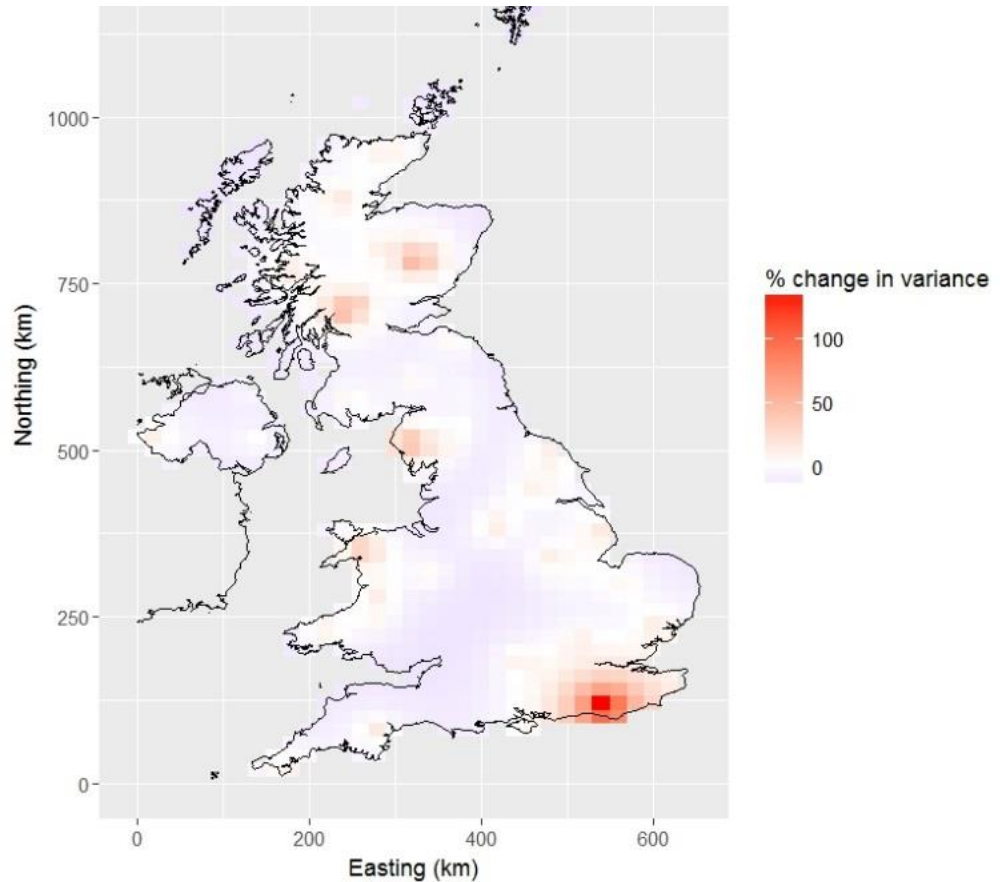
SPATIAL & SPATIO-
TEMPORAL

Generalise: Impact of sample size on spatial interpolation

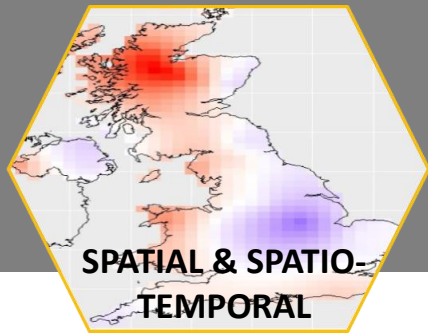
Change in estimated ammonium concentration (mean)



Change in variance of ammonium predictions





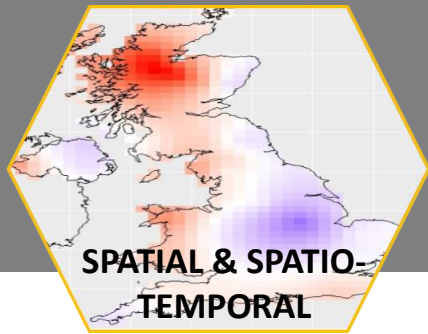


ESTIMATING THE MODEL

- Data - how are they distributed; do we have covariates?
- Mean - constant; trends
- Covariance model
- Priors – which are fixed, which are random, which distribution?

Krige.bayes function
model.control

prior.control



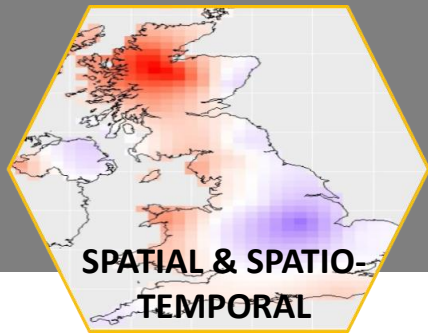
ESTIMATING THE MODEL

- Data - how are they distributed; do we have covariates?
- Mean - constant; trends
- Covariance model
- Priors – which are fixed, which are random, which distribution?
- Posterior distribution for each parameter
- Predictive distribution for each location
 - From which we can map the mean, variance etc...

Krige.bayes function
model.control

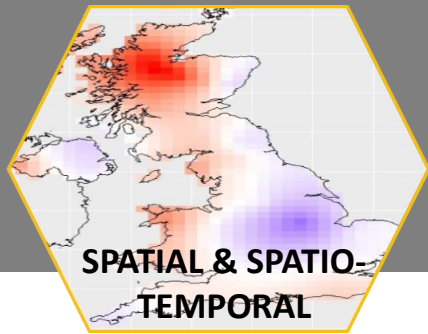
prior.control

output.control



SUMMARY

- In environmental data, things most often positively correlate in space
- Reflecting underlying processes and (un)measured data
- This may be of interest or a dependence for which we need to account
- Computational power is allowing us to do this in more sophisticated ways
- In Bayesian geostatistics the variogram parameters are random (defined by a probability distribution), not fixed (single value)
- As such we can propagate and more thoroughly estimate uncertainty, which may inform model selection as well as inference

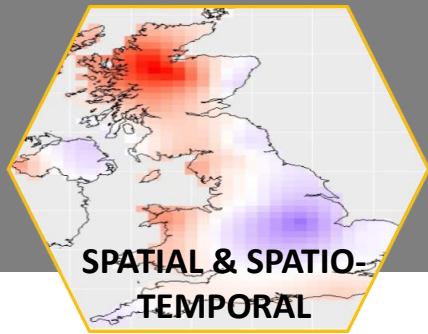


RESOURCES

- Banerjee & Fuentes (2012) Bayesian Modelling for Large Spatial Datasets. *Wiley Interdiscip Rev Comput Stat.* 4(1): 59–66.
- Oliver & Webster (2014) A tutorial guide to geostatistics: Computing and modelling variograms and kriging. *Catena* 113: 56-69
- Ribeiro et al. (2003) Geostatistical software – geoR and geoRglm. DSC 2003 Working Papers
- Ribeiro & Diggle Technical Report ST-99-08: Bayesian inference in Gaussian model-based geostatistics

See also references in Spatio-temporal slides

Supplementary slides



GEOSTATISTICS

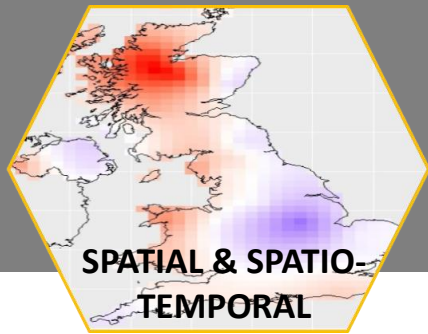
The mean function – the expected value at location \mathbf{s} from the distribution of all possible values generated from stochastic process $Z(\mathbf{s})$

$$\mu_Z(\mathbf{s}) = \mathbb{E}[Z(\mathbf{s})] \quad \mathbf{s} \in D$$

When $Z(\mathbf{s})$ is a continuous random variable:

$$\mu_Z(\mathbf{s}) = \mathbb{E}[Z(\mathbf{s})] = \int_{-\infty}^{\infty} z f_{Z(\mathbf{s})}(z) dz$$

Where $f_{Z(\mathbf{s})}$ is the probability density function (pdf) for $Z(\mathbf{s})$



GEOSTATISTICS

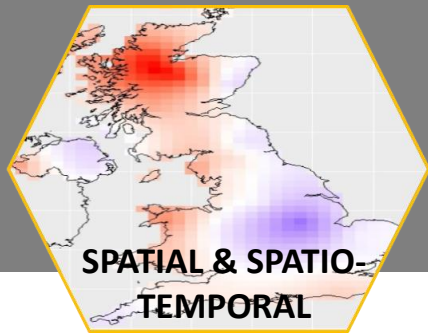
The covariance function – covariance measures strength of linear dependence between two random variables ($Z(s)$ and $Z(t)$)

$$C_Z(\mathbf{s}, \mathbf{t}) = \text{Cov}[Z(\mathbf{s}), Z(\mathbf{t})]$$

$$C_Z(\mathbf{s}, \mathbf{t}) = \mathbb{E}[(Z(\mathbf{s}) - \mu_Z(\mathbf{s}))(Z(\mathbf{t}) - \mu_Z(\mathbf{t}))]$$

The variance function of $Z(s)$ is the special case of the covariance $s=t$, giving:

$$\begin{aligned}\text{Var}[Z(\mathbf{s})] &= \text{Cov}[Z(\mathbf{s}), Z(\mathbf{s})] \\ &= \mathbb{E}[(Z(\mathbf{s}) - \mu_Z(\mathbf{s}))^2] \\ &= \varepsilon_Z^2(\mathbf{s})\end{aligned}$$



GEOSTATISTICS

The correlation function – the strength of association between two random variables ($Z(s)$ and $Z(t)$) is simply a scaled version of the covariance function:

$$\begin{aligned}\rho_z(\mathbf{s}, \mathbf{t}) &= \text{Corr}[Z(\mathbf{s}), Z(\mathbf{t})] \\ &= \frac{C_z(\mathbf{s}, \mathbf{t})}{\sqrt{C_z(\mathbf{s}, \mathbf{s})C_z(\mathbf{t}, \mathbf{t})}}\end{aligned}$$