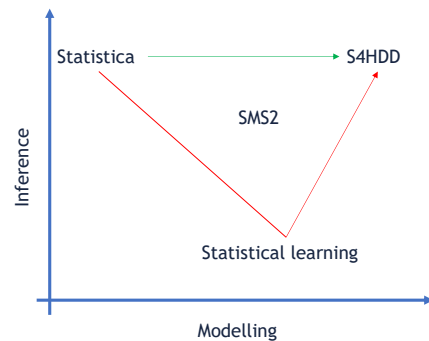# Statistics for High Dimensional Data (S4HDD) and CompStat Lab
## a.a. 2023/2024 (2nd edition)

Prof. Francesco Finazzi francesco.finazzi@unibg.it
Prof. Alessandro Fassò alessandro.fasso@unibg.it

UNIVERSITÀ DEGLI STUDI DI BERGAMO | Dipartimento di Ingegneria Gestionale, dell'Informazione e della Produzione

---

## Course structure

- **9 CFUs**
  - 3 CFUs on time series modelling (ex SMS2) (Prof. Fassò)
  - 6 CFUs on space and space-time data modelling (Prof. Finazzi)
  - Space-time modelling will be addressed late in the course
- **Lectures**
  - Theory
  - Coding (MATLAB)
  - Analysis of data sets
  - Scientific article discussion

UNIVERSITÀ DEGLI STUDI DI BERGAMO | Dipartimento di Ingegneria Gestionale, dell'Informazione e della Produzione

---

## Roadmap



UNIVERSITÀ DEGLI STUDI DI BERGAMO | Dipartimento di Ingegneria Gestionale, dell'Informazione e della Produzione

---

## Why this course

- To provide advanced statistical models/methods/tools that enable the extraction of useful information from temporal, spatial and spatio-temporal data sets (possibly large)...
- ...in order to make decisions...
- ...following a statistical inference approach...
- ...thus, knowing the risk to make the wrong decision.

UNIVERSITÀ DEGLI STUDI DI BERGAMO | Dipartimento di Ingegneria Gestionale, dell'Informazione e della Produzione

## Prerequisites

- Statistical inference
  - Distributions
  - Estimators
  - Confidence intervals
  - Hypothesis testing
- Simple and multiple regression
- Validation and cross-validation techniques
- Monte Carlo and bootstrap techniques
- Calcolo numeric (6 CFU) Prof. Maggioni

## Textbooks

- Time series analysis and its applications: with R examples, Robert H. Shumway, David S. Stoffer
- Model-based Geostatistics, Peter J. Diggle, Paulo J. Ribeiro
- Statistics for spatio-temporal data, Noel A. Cressie, Christopher K. Wikle
- Spatio-temporal statistics with R, Christopher K. Wikle, Andrew Zammit-Mangion, Noel Cressie

## Student evaluation

- Teamwork - Groups of 1 or 2 students(?)
- You choose the data set to analyze
- The analysis is based on the models/methods seen during the course
- You can use R, MATLAB or Python
- The evaluation will be based on a report
- Each group will present in front of the class (in English)
- More details during the course…

## Software

- MATLAB, R (and Python) are the most common environments used by researchers in space-time data analysis.
- Software for space-time data analysis is often open source (and not fully tested - use at your own risk!)
- In this course:
  - R is used for time series data analysis
  - MATLAB is used for space and space-time data analysis
- R and MATLAB are not fully interchangeable when it comes to "complex" statistical models and methods
- D-STEM (v2) will be the MATLAB software for space-time data analysis (download from Journal of Statistical Software)

## Examples of spatio-temporal data sets

- COVID19 pandemic data sets (hospitalizations and casualties by country over time)
- Climate change data sets (satellite measurements of temperature, humidity, etc.)
- Air quality data sets ($PM_{10}$, $PM_{2.5}$, NOx etc. observed at monitoring stations)
- Smartphone-based real time detection of earthquakes (www.sismo.app)
- Internet data analytics (app/web visits by country/region over time)
- People mobility data (www.facebook.com/covid19mobility)
- Digital image/video analysis (tracking, object recognition, etc.)
- Ecology (point processes)
- In general, addressing interesting societal problems requires the analysis of space-time data sets

## Why modelling spatio-temporal data

- For understanding the underlying data generation process
- For temporal, spatial and spatio-temporal prediction
    - Methods differ for continuous/discrete space and/or time
- For emulating complex phenomena (what is a model?)
- For making decisions

## Modelling steps

- Model proposal (which equation, which covariates, which random variables, which relations between variables)
- Model estimation (maximum likelihood, expectation-maximization algorithm)
- Model validation (k-fold cross-validation, leave-one-out validation)
- Model adoption (offline or online)
- Model update/improvement (when/if new data become available)

## A first spatial model

$$y(s) = x(s)'\beta + \varepsilon(s) \quad (1)$$

- $y(s)$ is the observation at generic spatial location $s \in R^2$ or $s \in S^2$
- $x(s)$ are spatial covariates (no missing data)
- $\beta$ is a vector of parameters
- $\varepsilon(s)$ is the random error at spatial location $s$ (e.g., $\varepsilon(s) \sim NID(0, \sigma_\varepsilon^2)$)
- $R^2$ is the 2D space (plane) while $S^2$ is the sphere embedded in $R^3$

## Model complexity

- When dealing with spatio–temporal data sets model complexity is often an issue
- Elements of model complexity
  - Univariate, bivariate, multivariate models
  - Number of covariates (p>n problem)
  - Number of spatial and/or temporal points (big n problem)
  - Dependences among variables
  - Space-time (non)separability
  - Latent variables/processes
  - Non Gaussianity

## How to estimate the spatial model?

- As usual we need data
- For us a data set has this form:

$$\{(x_1(s_1), y_1(s_1)), \dots, (x_n(s_n), y_n(s_n))\}$$

$$x_i = (x_{i1}, \dots, x_{ip})', i = 1, \dots, n$$

## Model residuals

$$e(s) = y(s) - x(s)'\widehat{\beta}$$
$$e(s) = y(s) - \hat{y}(s)$$

- We check if residuals $e(s)$ are IID which means:
  - Are residuls spatially uncorrelated?
  - Is the residul variance constant in space? (Homoscedasticity)

## Towards more complex models

- In general, residuals $e(s)$ are spatially correlated
- Why? Because $x(s)'\beta$ cannot entirely capture the data variability
- This is similar to the case of serially/temporally correlated residuals in classic regression models ($x'\beta$)
- When this happens we might
  - Add covariates
  - Add transformations of covariates (polynomials)
  - Add interactions between covariates
- This may not be enough to have IID residuals

## Towards more complex models

- What happens if we ignore the spatial correlation of residuals?
  - The model fitting capability is lower than it should be
  - Variance of estimators is wrong
  - Confidence intervals might have significance levels lower than their nominal levels
  - Spatial predictions are poor (especially in model validation)

## Spatial model with latent variable

$$y(s) = x(s)'\beta + \alpha w(s) + \varepsilon(s) \quad (2)$$

- $y(s)$, $x(s)'\beta$ and $\varepsilon(s)$ are the usual terms of a regression model (see model 1 in previous slides)
- $w(s)$ is a latent (non observed) random variable spatially correlated with unitary variance
- $corr(w(s), w(s')) = \rho(s, s'; \theta)$
- $\rho(s, s'; \theta)$ is a bivariate function (on $s$ and $s'$) with unknown parameter vector $\theta$
- $\alpha$ is a scale parameter to be estimated

## Spatial model with latent variable

$$y(s) = x(s)'\beta + \alpha w(s) + \varepsilon(s) \quad (2)$$

- $w(s) \perp \varepsilon(s)$, or $cov(w(s), \varepsilon(s)) = 0$
- Conditionally on $w(s)$, observed data $y_i(s_i)$ are realizations of mutually independent random variables $Y_i$ with conditional mean $E(Y_i|w) = x_i(s_i)'\beta + \alpha w(s_i)$ and conditional variance $\sigma_\varepsilon^2$
- What is the difference between $s_i$ and $s$?

## Spatial model with latent variable

$$y(s) = x(s)'\beta + \alpha w(s) + \varepsilon(s) \quad (2)$$

- $w(s)$ is spatially varying (in general each location $s$ has a different $w(s)$ value)
- Is $\varepsilon(s) \equiv 0$ (and $e(s) = 0$) if $w(s)$ can be choosen to explain $y(s)$ at each spatial location?
- How overfitting is avoided?
- What is the role of spatial correlation?

## Spatial correlation function

- Correlation function $\rho(s, s'; \theta)$ can be any postive-definite function
- This condition imposes that the linear combination $\sum_{i=1}^{m} a_i w(s_i)$ has positive variance
- $\theta$ is a vactor of unknown parameters (can be of any size)

## Examples of spatial correlation functions

- $\rho(s, s'; \theta)$ depends on spatial coordinates $s$ and $s'$
- Such a function describe the most «flexible» spatial correlation
- In 1D, the correlation between $w(0)$ and $w(1)$ would be different from the correlation between $w(1)$ and $w(2)$
- In most cases, $\rho(s, s'; \theta) = \rho(\|s - s'\|; \theta) = \rho(u; \theta)$ where $\| \ \|$ is the distance (Euclidean or geodetic) between $s$ and $s'$.
- In this case the spatial correlation only depends on the distance $u$ and not on the coordinates.

## Exponential spatial correlation function

- Exponential spatial correlation function

$$\rho(u; \theta) = exp(-u/\theta)$$

- $\theta > 0$ is a scalar value
- Simple but not flexible
- Not always suitable to model spatial correlation

## Matérn spatial correlation function

- Matérn correlation function:

$$\rho(u; \theta) = \{2^{\kappa-1}\Gamma(\kappa)\}^{-1}(u/\phi)K_\kappa(u/\phi)$$

- $K_\kappa$ is the modified Bessel function of the second kind of order $\kappa$
- The parameter $\kappa > 0$ is called the order of the Matérn and determines the differentiability of $w(s)$ (more details later)
- The parameter $\phi > 0$ determines the rate at which the correlation decays to zero with increasing $u$

## Gaussian process

- $w(s)$ is not directly observed but can be inferred from the observed data set $\{(x_1(s_1), y_1(s_1)), \dots, (x_n(s_n), y_n(s_n))\}$
- First, the distribution of $w(s)$ must be specified
- Formally, $w(s)$ is a zero mean Gaussian process (GP) with unitary variance and spatial correlation function $\rho(u; \boldsymbol{\theta})$
- $w(s)$ is a gaussian process if the joint distribution of $w(s_1), \dots, w(s_n)$ is multivariate normal (see SMS2) for any integer $n$ and any set of locations $s_1, \dots, s_n$

## Simulations

- How to simulate a GP with a given spatial correlation function in MATLAB?
- Simulations are useful to understand if a given correlation function is suitable to model the observed data set
- A GP is continuous in space. This means that:
    - It can be simulated on a regular grid or
    - On a irregular grid
    - For any given number of spatial locations
- When the GP is simulated on a regular grid, the simulated values referes to the centres of the pixels

## Statistics for High Dimensional Data (and CompStat Lab)
### a.a. 2023/2024 (2nd edition)

Prof. Francesco Finazzi francesco.finazzi@unibg.it
Prof. Alessandro Fassò alessandro.fasso@unibg.it

## Lesson 2

## Spatial model with latent variable

$$y(s) = x(s)'\beta + \alpha w(s) + \varepsilon(s) \quad (2)$$

- $w(s) \sim GP\big(0, \rho(\|s - s'\|; \theta)\big)$
- $\rho(\|s - s'\|; \theta) = corr\big(w(s), w(s')\big)$
- $\varepsilon(s) \sim N(0, \sigma_\varepsilon^2)$
- The unknown parameter set is $\Psi = \{\beta, \alpha, \sigma_\varepsilon^2, \theta\}$
- How to estimate $\Psi$ from data?

## Maximum likelihook estimate

- We rely on MLE to estimate the model parameter vector
- Because MLE has good properties and…
- Because we can use the EM algorithm (if needed)
- Likelihood function of (2) is:

$$L(\Psi; y, w, X) = L(\Psi; y|w, X)L(\Psi; w)$$

- $y = (y_1, \dots, y_n)'$ is the vector of observations at $n$ spatial locations
- $w = (w_1, \dots, w_n)'$ is the vector of latent variables at $n$ spatial locations
- $X$ is the $n \times p$ design matrix

## Likelihood decomposition

$$
\begin{aligned}
L(\Psi; y, w, X) &= L(\Psi; y|w, X)L(\Psi; w) \\
&= L(\beta, \alpha, \sigma_\varepsilon^2; y|w, X)L(\theta; w)
\end{aligned}
$$

- Each likelihood term depends on a subset of $\Psi$.
- $L(\Psi; y, w, X)$ is the complete-data likelihood (which assumes $w$ to be known)
- $L(\beta, \alpha, \sigma_\varepsilon^2; y|w, X)$ and $L(\theta; w)$ are densities of $n-$variate normal distributions

## Log-likelihood function

- As usual we prefer to work with $\log(L_\Psi)$
- $-2\log L_\Psi$ is given by:

$$\log|\Sigma_\varepsilon| + e'\Sigma_\varepsilon^{-1}e + \log|\Sigma_w| + w'\Sigma_w^{-1}w$$

where

- $e = y - X\beta - \alpha w$
- $\Sigma_\varepsilon = \sigma_\varepsilon^2 I_n$, with $I_n$ the identity matrix of dimension $n$
- $\Sigma_w$ is the $n \times n$ correlation matrix (e.g., $exp(-D/\theta)$, with $D$ the distance matrix)

## ML estimate

- MLE is given by

$$\widehat{\Psi} = argmin_{\boldsymbol{\beta},\alpha,\sigma_\varepsilon^2,\boldsymbol{\theta}} \quad log|\Sigma_\varepsilon| + \boldsymbol{e}'\Sigma_\varepsilon^{-1}\boldsymbol{e} + log\ |\Sigma_w| + \boldsymbol{w}'\Sigma_w^{-1}\boldsymbol{w}$$

- Argmin because we are considering $-2\log(L_\Psi)$
- Unfortunately minimizing $-2\log(L_\Psi)$ is not feasible, plus $\boldsymbol{w}$ is latent and not observed
- We must rely on the EM algorithm

## EM algorithm

- The EM is an iterative algorithm for MLE
- First iteration starts with initial values $\widehat{\Psi}^{\langle 0 \rangle}$ (usually given by OLS and method of moments)

- E-step
$$Q\big(\Psi, \widehat{\Psi}^{\langle m \rangle}\big) = E_{\widehat{\Psi}^{\langle m \rangle}}(-2\mathrm{logL}(\Psi; \boldsymbol{y}, \boldsymbol{w}, \boldsymbol{X})|\boldsymbol{y})$$

- M-step
$$\widehat{\Psi}^{\langle m+1 \rangle} = argmax_\Psi\ Q\big(\Psi, \widehat{\Psi}^{\langle m \rangle}\big)$$

## EM algorithm, E-step

- E-step

$$E_{\widehat{\Psi}^{\langle m \rangle}}(-2\mathrm{logL}(\Psi; \boldsymbol{y}, \boldsymbol{w}, \boldsymbol{X})|\boldsymbol{y})$$
$$= tr\big[\Sigma_\varepsilon^{-1}\big(E(\boldsymbol{e}|\boldsymbol{y})E(\boldsymbol{e}|\boldsymbol{y})' + Var(\boldsymbol{e}|\boldsymbol{y})\big)\big]$$
$$+ tr\big[\Sigma_w^{-1}\big(E(\boldsymbol{w}|\boldsymbol{y})E(\boldsymbol{w}|\boldsymbol{y})' + Var(\boldsymbol{w}|\boldsymbol{y})\big)\big]$$

- $E(\boldsymbol{e}|\boldsymbol{y}) = \boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta} - \alpha E(\boldsymbol{w}|\boldsymbol{y})$
- $Var(\boldsymbol{e}|\boldsymbol{y}) = Var(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta} - \alpha\boldsymbol{w}|\boldsymbol{y}) = \alpha^2 Var(\boldsymbol{w}|\boldsymbol{y})$

## EM algorithm, E-step

- $E(\boldsymbol{w}|\boldsymbol{y}) = Cov(\boldsymbol{w}, \boldsymbol{y})Var(\boldsymbol{y})^{-1}[\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}]$ (see multivariate normal)
- $Var(\boldsymbol{w}|\boldsymbol{y}) = \Sigma_w - Cov(\boldsymbol{w}, \boldsymbol{y})Var(\boldsymbol{y})^{-1}Cov(\boldsymbol{w}, \boldsymbol{y})'$
- $Var(\boldsymbol{y}) = Var(\boldsymbol{X}\boldsymbol{\beta} + \alpha\boldsymbol{w} + \boldsymbol{\varepsilon}) = Var(\alpha\boldsymbol{w} + \boldsymbol{\varepsilon}) = \alpha^2 Var(\boldsymbol{w}) + Var(\boldsymbol{\varepsilon}) + 2Cov(\boldsymbol{w}, \boldsymbol{\varepsilon})$
- $Var(\boldsymbol{w}) = \Sigma_w$
- $Var(\boldsymbol{\varepsilon}) = \Sigma_\varepsilon = \sigma_\varepsilon^2 \boldsymbol{I}_n$
- $2Cov(\boldsymbol{w}, \boldsymbol{\varepsilon}) = \boldsymbol{0}$ (from model assumptions)
- $Cov(\boldsymbol{w}, \boldsymbol{y}) = Cov(\boldsymbol{w}, \boldsymbol{X}\boldsymbol{\beta} + \alpha\boldsymbol{w} + \boldsymbol{\varepsilon}) = Cov(\boldsymbol{w}, \alpha\boldsymbol{w}) = \alpha Cov(\boldsymbol{w}, \boldsymbol{w}) = \alpha Var(\boldsymbol{w}) = \alpha\Sigma_w$

## EM algorithm, M-step

$$\widehat{\Psi}^{\langle m+1 \rangle} = argmax_\Psi \, Q\big(\Psi, \widehat{\Psi}^{\langle m \rangle}\big)$$

$$\frac{dQ\big(\Psi, \widehat{\Psi}^{\langle m \rangle}\big)}{d\Psi} = 0$$

$$\alpha^{\langle m+1 \rangle} = \frac{tr\big[(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}^{\langle m \rangle}) E(\boldsymbol{w}|\boldsymbol{y})'\big]}{tr[E(\boldsymbol{w}|\boldsymbol{y}) E(\boldsymbol{w}|\boldsymbol{y})' + Var(\boldsymbol{w}|\boldsymbol{y})]}$$

$$\boldsymbol{\beta}^{\langle m+1 \rangle} = (\boldsymbol{X}'\boldsymbol{X})^{-1} \left[ \boldsymbol{X}' \left( \boldsymbol{y} - \alpha^{\langle m+1 \rangle} E(\boldsymbol{w}|\boldsymbol{y}) \right) \right]$$

---

## EM algorithm, M-step

$$\sigma_\varepsilon^{2\,\langle m+1 \rangle} = \frac{1}{n} tr[E(\boldsymbol{e}|\boldsymbol{y}) E(\boldsymbol{e}|\boldsymbol{y})' + Var(\boldsymbol{e}|\boldsymbol{y})]$$

$$\boldsymbol{\theta}^{\langle m+1 \rangle} = argmin_{\boldsymbol{\theta}} \; log \, |\Sigma_w^{-1}(\boldsymbol{\theta})| + tr[\Sigma_w^{-1}(\boldsymbol{\theta})(\widehat{\boldsymbol{w}}\widehat{\boldsymbol{w}}')]$$

Where $\widehat{\boldsymbol{w}} = E_{\boldsymbol{\theta}^{\langle m \rangle}}(\boldsymbol{w}|\boldsymbol{y}) = E(\boldsymbol{w}|\boldsymbol{y})$

---

# Statistics for High Dimensional Data (and CompStat Lab)
## a.a. 2023/2024 (2nd edition)

Prof. Francesco Finazzi francesco.finazzi@unibg.it
Prof. Alessandro Fassò alessandro.fasso@unibg.it

---

# Lesson 3

## Spatial model with latent variable

$$y(s) = x(s)'\boldsymbol{\beta} + \alpha w(s) + \varepsilon(s) \quad (2)$$

- $w(s) \sim GP\big(0, \rho(\|s - s'\|; \boldsymbol{\theta})\big)$
- $\rho(\|s - s'\|; \boldsymbol{\theta}) = corr(w(s), w(s'))$
- $\varepsilon(s) \sim N(0, \sigma_\varepsilon^2)$
- The unknown parameter set is $\Psi = \{\boldsymbol{\beta}, \alpha, \sigma_\varepsilon^2, \boldsymbol{\theta}\}$
- $\Psi$ is estimated using the EM algorithm

---

## Prediction using model

- Once estimated, the model is used for prediction:

$$\hat{y}(s_i) = x(s_i)'\widehat{\boldsymbol{\beta}} + \hat{\alpha}\widehat{w}(s_i)$$

- The above formula gives the prediction at spatial locations $s_i, i = 1, \dots, n$
- $\widehat{w}(s_i) = E(w(s_i)|Y)$
- How to estimate $y$ at a generic spatial location $s$?
- (First of all we need $x(s)$)

---

## Spatial prediction

- First of all we need $x(s)$ (spatial covariates at $s$)
- Remember that $E(w|y) = Cov(w, y)Var(y)^{-1}[y - X\beta]$
- Similarly $\widehat{w}(s) = E(w(s)|y) = Cov(w(s), y)Var(y)^{-1}[y - X\beta]$
- In practice the prediction $\widehat{w}(s)$ depends on the vector $y$ of all observations
- More in details:
  $Cov(w(s), y) = Cov(w(s), X\beta + \alpha w + \varepsilon) = Cov(w(s), \alpha w) = \alpha Cov(w(s), w)$
- $Cov(w(s), w)$ is a $1 \times n$ vector with elements $Cov\big(w(s), w(s_i)\big) = exp\left(-\frac{\|s - s_i\|}{\theta}\right)$.
- $Var(y)^{-1}[y - X\beta]$ are the same seen for the EM algorithm and do not depend on $w(s)$.

---

## Spatial prediction

- When prediction is done for multiple spatial locations $S = \{s_1, \dots, s_M\}$:

$$\widehat{w}(S) = E(w(S)|y) = Cov(w(S), y)Var(y)^{-1}[y - X\beta]$$

- $Cov(w(S), \alpha w) = \alpha Cov(w(S), w)$
- $Cov(w(S), w)$ is a $M \times n$ matrix
- $Var(y)^{-1}[y - X\beta]$ is computed only one time!
- $Var(w(S)|y)$ is the prediction uncertainty

- How to make spatial prediction with D-STEM?

## Multivariate models

- It is not uncommon to jointly model multiple variables (e.g., multiple pollutants, multiple met. variables, etc.)
- The spatial model becomes multivariate

$$y(s) = X(s)'\beta + \alpha w(s) + \varepsilon(s) \quad (2)$$

- $y(s)$ is a $p \times 1$ vector
- $w(s) \sim GP_p\big(\mathbf{0}, V\rho(\|s - s'\|; \theta)\big)$ is a p-variate Gaussian random process
- $V$ is a $p \times p$ correlation matrix
- $\varepsilon(s) \sim N_p(\mathbf{0}, \Sigma_\varepsilon^2)$ is a p-variate Normal random variable with $\Sigma_\varepsilon^2$ diagonal
- The unknown parameter set is $\Psi = \{\beta, \alpha, \Sigma_\varepsilon^2, \theta, V\}$

## Why a multivariate model?

- Why not fitting a model for each variable?
- If two or more variables are correlated, spatial prediction can benefit from this correlation
- Especially if one variables is observed at few spatial locations w.r.t. the other variables
- But computing time is higher (matrices are $pn \times pn$ if all variables are observed at $n$ locations)

## Bivariate model

- A bivariate model is a (simple) special case of the multivariate model

$$\begin{bmatrix} y_1(s) \\ y_2(s) \end{bmatrix} = \begin{bmatrix} x_1(s)' & \mathbf{0} \\ \mathbf{0} & x_2(s)' \end{bmatrix} \begin{bmatrix} \beta_1(s) \\ \beta_2(s) \end{bmatrix} + \begin{bmatrix} \alpha_1 w_1(s) \\ \alpha_2 w_2(s) \end{bmatrix} + \begin{bmatrix} \varepsilon_1(s) \\ \varepsilon_2(s) \end{bmatrix}$$

- $x_1(s)'$ and $x_2(s)'$ can have different lengths
- $V = \begin{bmatrix} 1 & corr(w_1(s), w_2(s)) \\ corr(w_2(s), w_1(s)) & 1 \end{bmatrix}$
- $\Sigma_\varepsilon^2 = \begin{bmatrix} \sigma_{1,\varepsilon}^2 & 0 \\ 0 & \sigma_{2,\varepsilon}^2 \end{bmatrix}$

## Linear model of coregionalization

- $V\rho(\|s - s'\|; \theta)$ is called linear coregionalization model
- $\theta$ is common to the $p$ variables!
- This could be a limit because all the variables are forced to share the same spatial correlation (same function and strenght)
- $V$ is a symmetric correlation matrix, only $p(p-1)/2$ elements of $V$ are estimated

## Data structure

- A multivatiate data set can be classified depending on spatial locations:
  - **Isotopic**: all $p$ variables are observed at the same $n$ spatial locations.
  - **Fully heterotipic**: the $p$ variables do not share a single spatial location.
  - **Partially heterotopic**: some of the $p$ variables are observed at a subset of the $n$ spatial locations. This is the most common case and it is the case handled by D-STEM.

# Statistics for High Dimensional Data (and CompStat Lab)
# a.a. 2023/2024 (2st edition)

bitvector,bitvector

Prof. Francesco Finazzi francesco.finazzi@unibg.it
Prof. Alessandro Fassò alessandro.fasso@unibg.it

# Lesson 4

## Spatio-temporal model: the dynamic coregionalization model (DCM)

$$
\begin{aligned}
y(\boldsymbol{s}, t) &= \boldsymbol{x_\beta}(\boldsymbol{s}, t)' \boldsymbol{\beta} + \boldsymbol{x_z}(\boldsymbol{s})' \boldsymbol{z}(t) + \alpha w(\boldsymbol{s}, t) + \varepsilon(\boldsymbol{s}, t) \\
\boldsymbol{z}(t) &= \boldsymbol{G}\boldsymbol{z}(t-1) + \boldsymbol{\eta}(t)
\end{aligned}
$$

- $\boldsymbol{x_\beta}(\boldsymbol{s}, t)$ and $\boldsymbol{x_z}(\boldsymbol{s})$ are vectors of covariates. Note that $\boldsymbol{x_z}(\boldsymbol{s})$ is time invariant.
- $w(\boldsymbol{s}, t) \sim GP\big(0, \rho(\|\boldsymbol{s} - \boldsymbol{s}'\|; \boldsymbol{\theta})\big)$ is correlated over space but IID over time
- $\boldsymbol{z}(t)$ is $q \times 1$ dimensional with Markovian dynamics
- $\boldsymbol{G}$ is a stable $q \times q$ transition matrix
- $\boldsymbol{\eta}(t) \sim N\big(\boldsymbol{0}, \boldsymbol{\Sigma}_\eta\big)$ is the innovation with $\boldsymbol{\Sigma}_\eta$ the variance–covariance matrix
- $\varepsilon(\boldsymbol{s}, t) \sim N(0, \sigma_\varepsilon^2)$ is the measurement error

## Why the Markovian dynamics?

- The Markovian dynamic helps to describe the temporal persistence which usually characterizes temporal phenomena.
- For instance: even if we stop air pollution emissions, pollutant concentration will not drop instantly to zero.
  - Emission would be a model covariate

## Data matrix

- $Y = (y_1, \dots, y_T)$ is the data matrix
- Each $y_t$ is the vector of spatial observation at time $t = 1, \dots, T$
- In each $y_t$, missing data are possible

- Covariates are in a data array $X = \{X_1, \dots, X_T\}$, each $X_t$ is a $n \times b$ matrix where $b$ is the number of covariates. $X$ cannot have missing data.

## Likelihood function

- The parameter vector is $\Psi = \{\boldsymbol{\beta}, \alpha, \sigma_\varepsilon^2, \boldsymbol{\theta}, \boldsymbol{G}, \boldsymbol{\Sigma}_\eta\}$
- Likelihood function is:

$$L(\Psi; Y, W, Z, X) = L(\Psi; Y|W, Z, X)L(\Psi; Z)L(\Psi; W)$$

- $W = (w_1, \dots, w_T)$
- $Z = (z_1, \dots, z_T)$

## Log-likelihood function

- $-2\log L_\Psi$ is given by:

$$T\log|\boldsymbol{\Sigma}_\varepsilon| + \sum_{t=1}^{T} e_t' \boldsymbol{\Sigma}_\varepsilon^{-1} e_t + T\log|\boldsymbol{\Sigma}_\eta| + \sum_{t=1}^{T}(z_t - Gz_{t-1})' \boldsymbol{\Sigma}_\eta^{-1}(z_t - Gz_{t-1}) + T\log|\boldsymbol{\Sigma}_w| + \sum_{t=1}^{T} w_t' \boldsymbol{\Sigma}_w^{-1} w_t$$

where
- $e_t = y_t - X_{\beta,t}\boldsymbol{\beta} - X_{z,t}z_t - \alpha w_t$
- $\boldsymbol{\Sigma}_\varepsilon = \sigma_\varepsilon^2 I_n$, with $I_n$ the identity matrix of dimension $n$
- $\boldsymbol{\Sigma}_w$ is the $n \times n$ correlation matrix (e.g., $exp(-D/\theta)$, with $D$ the distance matrix)

## Model estimation – EM algorithm

- Model estimation is (again) based on the EM algorithm
- $E(w(\boldsymbol{s},t)|\boldsymbol{Y})$ and $Var(w(\boldsymbol{s},t)|\boldsymbol{Y})$ are given by the same formulas of $E(w(\boldsymbol{s})|\boldsymbol{Y})$ and $Var(w(\boldsymbol{s})|\boldsymbol{Y})$ (because $w(\boldsymbol{s},t)$ are IID over time)
- $E(\boldsymbol{z}(t)|\boldsymbol{Y})$ and $Var(\boldsymbol{z}(t)|\boldsymbol{Y})$ are given by the Kalman smoother
- However, $Cov(\boldsymbol{z}(t),w(\boldsymbol{s},t)|\boldsymbol{Y}) \neq \boldsymbol{0}$. E-step and M-step are more complicated than the spatial case.

---

# Statistics for High Dimensional Data (and CompStat Lab)
# a.a. 2023/2024 (2nd edition)

Prof. Francesco Finazzi  francesco.finazzi@unibg.it
Prof. Alessandro Fassò  alessandro.fasso@unibg.it

---

# Lesson 5

---

## Spatio-temporal model: the (univariate) hidden dynamic geostatistical model (HDGM)

$$
\begin{aligned}
y(\boldsymbol{s},t) &= \boldsymbol{x_\beta}(\boldsymbol{s},t)'\boldsymbol{\beta} + az(\boldsymbol{s},t) + \varepsilon(\boldsymbol{s},t) \\
z(\boldsymbol{s},t) &= gz(\boldsymbol{s},t-1) + \eta(\boldsymbol{s},t)
\end{aligned}
$$

- $\eta(\boldsymbol{s},t) \sim GP\big(0, \rho(\|\boldsymbol{s}-\boldsymbol{s}'\|;\boldsymbol{\theta})\big)$ is correlated over space but IID over time
- $z(\boldsymbol{s},t)$ is scalar and has Markovian dynamic
- $a$ is a scale coefficient ($v$ in D-STEM)
- $g$ is the transition coefficient
- $\varepsilon(\boldsymbol{s},t) \sim N(0,\sigma_\varepsilon^2)$ is the measurement error
- The model parameter set is $\Psi = \{\boldsymbol{\beta}, a, \sigma_\varepsilon^2, \boldsymbol{\theta}, g\}$

- Which are the main differences with the DCM?
- Which model is better?

## Spatio-temporal model: the (multivariate) hidden dynamic geostatistical model (HDGM)

$$
\begin{aligned}
\boldsymbol{y}(\boldsymbol{s}, t) &= \boldsymbol{X_\beta}(\boldsymbol{s}, t)'\boldsymbol{\beta} + \boldsymbol{z}(\boldsymbol{s}, t) + \boldsymbol{\varepsilon}(\boldsymbol{s}, t) \\
\boldsymbol{z}(\boldsymbol{s}, t) &= \boldsymbol{G}\boldsymbol{z}(\boldsymbol{s}, t-1) + \boldsymbol{\eta}(\boldsymbol{s}, t)
\end{aligned}
$$

- $\boldsymbol{y}(\boldsymbol{s}, t)$ and $\boldsymbol{z}(\boldsymbol{s}, t)$ are $p \times 1$ vectors
- $\boldsymbol{\eta}(\boldsymbol{s}, t) \sim GP_p\big(\boldsymbol{0}, \boldsymbol{V}p(\|\boldsymbol{s}-\boldsymbol{s}'\|; \boldsymbol{\theta})\big)$ is a p-variate Gaussian random process
- $\boldsymbol{V}$ is a variance-covariance matrix (in Calculli et al. $\boldsymbol{V}$ is a correlation matrix and there is the scaling matrix $\boldsymbol{A}$)
- $\boldsymbol{z}(\boldsymbol{s}, t)$ has Markovian dynamic
- $\boldsymbol{G}$ is a diagonal $p \times p$ transition matrix
- $\boldsymbol{\varepsilon}(\boldsymbol{s}, t) \sim N_p(\boldsymbol{0}, \boldsymbol{\Sigma_\varepsilon})$ is a p-variate Normal random variable with $\boldsymbol{\Sigma_\varepsilon}$ diagonal
- The model parameter set is $\Psi = \{\boldsymbol{\beta}, \boldsymbol{\Sigma_\varepsilon}, \boldsymbol{\theta}, \boldsymbol{V}, \boldsymbol{G}\}$

---

## Model estimation

- The HDGM is estimated similarly to the DCM
- But we only have the $z(\boldsymbol{s}, t)$ latent variable which is estimated in the E-step by the Kalman smoother.
- Spatial prediction is also done by the Kalman smoother assuming that y is not observed at the spatial prediction locations (it is added as NaN in the y vector).

---

# Statistics for High Dimensional Data (and CompStat Lab)
# a.a. 2023/2024 (2nd edition)

Prof. Francesco Finazzi francesco.finazzi@unibg.it

Prof. Alessandro Fassò alessandro.fasso@unibg.it

UNIVERSITÀ DEGLI STUDI DI BERGAMO | Dipartimento di Ingegneria Gestionale, dell'Informazione e della Produzione

# Lesson 6

## Towards spatio-temporal functional models

- In many cases, data are observed at high frequency/resolution in at least one dimension (spatial or temporal)
- It usually happens with the temporal dimension
  - For instance, pollutant concentrations observed hourly or every 15 minutes
  - In general, high frequency observations (100, 1000, 10.000 per day)
- In space it is less common because a high resolution sampling is usually very expensive but...
  - In 3D space, one dimension may be sampled at higher resolution than the others

## Towards spatio-temporal functional models

- Which are the problems with high frequency/resolution data?
  - Usually the original data set is very large and so the computational burden (of classic spatio-temporal models)
  - Data may be collected asynchronously over time (e.g., different monitoring stations may have different clocks)
  - Data may have large gaps over time (how does the Kalman Smoother perform in this case?)
  - Temporal correlation is usually very high (the Markovian model may explain the data but it is not very useful for prediction)

## Functional data analysis (FDA)

- In FDA, the object of the statistical inference is a continuous function rather than scalar/vector values
- For instance, the temperature measured by a sensor over the 24h of the day can be described by a (smooth?) function
  - Independently of how many observations we take
  - Independently of where in time these observations are taken
- Which function or class of functions should we use?
  - The function should describe the «global» data pattern
  - In a way, the function filters out the data noisy
  - The researcher should be able to control the function smoothness

## Splines

- Spline is a class of functions which allows to easily control the function smoothness
  - By selecting the proper basis functions
  - By selecting the proper knots
- B-spline basis are useful to describe non-periodic functions
  - Knots can be placed ad-hoc along the function domain (more knots where the function should change more rapidly)
- Fourier basis can describe periodic functions
  - Smoothness is controlled by the number of basis

## How to describe functional data in a space-time model?

- We now want to model the generic observation $y(\boldsymbol{s}, t, h)$
  - $\boldsymbol{s}$ and $t$ are the usual spatial and temporal indexes
  - $h \in \mathbb{R}$ is the «functional» dimension (spatial or temporal)
- Examples
  - $h$ could describe the continuous time within the day while $t$ is the index of days
  - $h$ could describe altitude in a 3D space while $\boldsymbol{s}$ describes the generic location across the globe

## The functional HDG model in D-STEM

- D-STEM implements the (univariate) functional version of the HDG model:

$$
\begin{aligned}
y(\boldsymbol{s}, t, h) &= \boldsymbol{x}(\boldsymbol{s}, t, h)' \boldsymbol{\beta}(h) + \boldsymbol{\phi}(h)' \boldsymbol{z}(\boldsymbol{s}, t) + \varepsilon(\boldsymbol{s}, t, h) \\
\boldsymbol{z}(\boldsymbol{s}, t) &= \boldsymbol{G} \boldsymbol{z}(\boldsymbol{s}, t-1) + \boldsymbol{\eta}(\boldsymbol{s}, t)
\end{aligned}
$$

- $\boldsymbol{\phi}(h)$ are the basis functions, $\boldsymbol{z}(\boldsymbol{s}, t)$ are the spline coefficients
- All details are in Wang et al. (2021) Journal statistical software