

# S4HDD

Silviu Filote

July 2023

## Contents

<b>1 Lesson 1</b>	<b>1</b>
1.1 Spatial model . . . . .	1
1.2 Models residuals . . . . .	1
1.3 Spatial model with latent variable . . . . .	2
1.4 Exponential spatial correlation function . . . . .	3
1.5 Matérn spatial correlation function . . . . .	3
1.6 Gaussian process (GP) . . . . .	3
1.7 Matlab exercises 2-3 . . . . .	4
1.8 Spatial correlation - Regular grid . . . . .	4
1.9 Spatial correlation - Irregular grid . . . . .	5
<b>2 Lesson 2</b>	<b>6</b>
2.1 Spatial model with latent variable . . . . .	6
2.2 Maximum likelihook estimate . . . . .	6
2.3 Likelihood decomposition . . . . .	7
2.4 Log-likelihood function . . . . .	7
2.5 ML estimate . . . . .	7
2.6 EM algorithm . . . . .	7
2.7 Variogram . . . . .	9
<b>3 Lesson 3</b>	<b>11</b>
3.1 Spatial model with latent variable . . . . .	11
3.2 Prediction for one site $s_i$ . . . . .	11
3.3 Spatial prediction for multiple sites $S$ . . . . .	12
3.4 Multivariate models . . . . .	12
3.5 Why a multivariate model? . . . . .	13
3.6 Bivariate model . . . . .	14
3.7 Linear model of coregionalization . . . . .	14
3.8 Data structure . . . . .	14
<b>4 Lesson 4</b>	<b>15</b>
4.1 Spatio-temporal model: the dynamic coregionalization model (DCM) . . . . .	15
4.2 Why the Markovian dynamics? . . . . .	15
4.3 Data matrix . . . . .	16
4.4 Likelihood function . . . . .	16
4.5 Log-likelihood function . . . . .	16
4.6 Model estimation – EM algorithm . . . . .	17
4.7 DCM alternative formulas . . . . .	17
<b>5 Lesson 5</b>	<b>18</b>
5.1 Spatio-temporal model: the (univariate) hidden dynamic geostatistical model (HDGM) . .	18
5.2 Spatio-temporal model: the (multivariate) hidden dynamic geostatistical model (HDGM)	18
5.3 Model estimation . . . . .	18

<b>6 Lesson 6</b>	<b>19</b>
6.1 Towards spatio-temporal functional models . . . . .	19
6.2 Functional data analysis (FDA) . . . . .	19
6.3 Basis Functions . . . . .	19
6.4 How to describe functional data in a space-time model? . . . . .	20
6.5 The functional HDG model in D-STEM . . . . .	20
<b>7 Apporfondimenti</b>	<b>21</b>
7.1 Kriging . . . . .	21
7.2 Gaussian process . . . . .	22
7.3 Apulia paper . . . . .	23
7.4 Project clarifications . . . . .	23
<b>8 Time series analysis</b>	<b>25</b>
8.1 Introduction: . . . . .	25
8.2 Autocorrelation (ACF) and partialcorrelation (PACF) function: . . . . .	25
8.3 Stationarity . . . . .	25
8.4 Unit Roots . . . . .	26
8.5 Dickey Fuller Test and Augmented Dickey Fuller Test . . . . .	27
8.6 White noise . . . . .	27
8.7 Backshift operator: lag operator . . . . .	27
8.8 Autoregressive Model . . . . .	28
8.9 Evaluate time series model . . . . .	28
<b>9 Moving average model</b>	<b>29</b>
9.1 Moving Average and ACF(Auto correlation function) . . . . .	29
9.2 Invertibility of Time Series: $MA(1) \Leftrightarrow AR(\infty)$ . . . . .	30
9.3 Invertibility of Time Series: $AR(1) \Leftrightarrow MA(\infty)$ . . . . .	30
9.4 ARMA model . . . . .	30
9.5 ARIMA model . . . . .	31
9.6 What is Seasonality . . . . .	31
<b>10 Roadmap time series</b>	<b>32</b>
<b>11 Machine learning models</b>	<b>33</b>
11.1 XGBoost . . . . .	33
11.2 Code: . . . . .	34
11.3 Understanding Facebook's Prophet . . . . .	34

# 1 Lesson 1

## 1.1 Spatial model

$$y(\mathbf{s}) = \mathbf{x}(\mathbf{s})' \cdot \boldsymbol{\beta} + \varepsilon(\mathbf{s})$$

Osservazioni:

- $y(\mathbf{s})$  is the observation at generic spatial location  $\mathbf{s} \in \mathbb{R}^2 = (\text{latitude}, \text{longitude})$  or  $\mathbf{s} \in \mathbb{S}^2$
- $\mathbf{x}(\mathbf{s})$  are generic spatial covariates,  $\mathbf{x}(\mathbf{s}_i)$  is the covariates in a specific location  $\mathbf{s}_i$
- $\boldsymbol{\beta}$  is a vector of parameters, it has to be estimated
- $\varepsilon(\mathbf{s})$  is the random error at spatial location  $\mathbf{s}$  (e.g.  $\varepsilon(\mathbf{s}) \sim NID(0, \sigma_\varepsilon^2)$ )
- $\mathbb{R}^2$  is the 2D space (plane) while  $\mathbb{S}^2$  is the sphere embedded in  $\mathbb{R}^3$
- $n$  locations  $N$  data  $\Rightarrow n = N$  (no multiple observations in the location, cause temporal)
- The data set has this form:

$$\begin{aligned} & \{(\mathbf{x}_1(\mathbf{s}_1), y_1(\mathbf{s}_1)), \dots, (\mathbf{x}_n(\mathbf{s}_n), y_n(\mathbf{s}_n))\} \\ & \mathbf{x}_i = (x_{i1}, \dots, x_{ip})', \quad i = 1, \dots, n \quad p = \text{covariates} \end{aligned}$$

## 1.2 Models residuals

$$\begin{aligned} e(t) &= \varepsilon(\mathbf{s}) = y(\mathbf{s}) - \mathbf{x}(\mathbf{s})' \cdot \hat{\boldsymbol{\beta}} \\ &= y(\mathbf{s}) - \hat{y}(\mathbf{s}) \end{aligned}$$

We check if residuals  $\varepsilon(\mathbf{s})$  are IID which means:

- **Are residuals spatially uncorrelated?** Most of the time the residuals are correlated due to the fact that we don't have all the information in the dataset, so we are not observing all the covariates. The model presented before is a **suboptimal model**.
- **Is the residual variance constant in space? (Homoscedasticity)** Heteroscedasticity means that we have information left inside the residuals, that's because our model didn't capture all the variability.

Observations:

- In general, residuals  $e(\mathbf{s})$  are **spatially correlated**, because  $\mathbf{x}(\mathbf{s})' \cdot \boldsymbol{\beta}$  cannot entirely capture the data variability. Probably this happens because measure all of the covariates can be expensive.
- When this happens we might (increase the complexity of the model):
  - Add covariates
  - Add transformations of covariates (polynomials)
  - Add interactions between covariates
$$\Rightarrow \text{This may not be enough to have IID residuals}$$

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  $\Rightarrow$  model assumptions is that  $\varepsilon(\mathbf{s})$  is IID and normally distributed
- Confidence intervals might have significance levels lower than their nominal levels
- Spatial predictions are poor (especially in model validation)
- We can't trust our model and our predictions

### 1.3 Spatial model with latent variable

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

Observations:

- $w(\mathbf{s})$  is a **latent (non observed) random variable** spatially correlated with unitary variance
- $w(\mathbf{s})$  is included as a component in the model to capture the spatial variation that is not explained by the observed covariates  $\mathbf{x}(\mathbf{s})$ . Including  $w(\mathbf{s})$  in the model, you acknowledge the presence of spatial correlation that cannot be explained by the covariates alone.
- $w(\mathbf{s})$  models the missing information and/or what the model doesn't capture
- $\varepsilon(\mathbf{s})$  models the measurement error + model error
- $w(\mathbf{s}) \perp \varepsilon(\mathbf{s})$ , or  $\text{cov}(w(\mathbf{s}), \varepsilon(\mathbf{s})) = 0$
- $\alpha$  is a scale parameter to be estimated  $\Rightarrow$  it is like a weight
- $\text{corr}(w(\mathbf{s}), w(\mathbf{s}')) = \rho(\mathbf{s}, \mathbf{s}', \boldsymbol{\theta})$
- $\rho(\mathbf{s}, \mathbf{s}', \boldsymbol{\theta})$  is a bivariate function (on  $\mathbf{s}$  and  $\mathbf{s}'$ ) with unknown parameter vector  $\boldsymbol{\theta}$
- Conditionally on  $w(\mathbf{s})$ , observed data  $y_i(\mathbf{s}_i)$  are realizations of mutually independent random variables  $Y_i$  with conditional mean  $E(Y_i|w) = \mathbf{x}(\mathbf{s})'\boldsymbol{\beta} + \alpha w(\mathbf{s})$  and conditional variance  $\sigma_\varepsilon^2$

#### Insights

- The latent variable  $w(\mathbf{s})$  is a random variable that **is spatially correlated**. In other words, the value of  $w(\mathbf{s})$  at one location  $\mathbf{s}$  is related to the values of  $w$  at nearby locations. This spatial correlation is quantified by a spatial correlation function, typically denoted as  $\rho(\mathbf{s}, \mathbf{s}', \boldsymbol{\theta})$ , where  $\boldsymbol{\theta}$  represents the parameters of this function.

$$\text{corr}(w(\mathbf{s}), w(\mathbf{s}')) = \rho(\mathbf{s}, \mathbf{s}', \boldsymbol{\theta})$$

- The spatial correlation function,  $\rho(\mathbf{s}, \mathbf{s}', \boldsymbol{\theta})$ , describes how the correlation between  $w(\mathbf{s})$  and  $w(\mathbf{s}')$  changes as a function of the spatial separation between locations  $\mathbf{s}$  and  $\mathbf{s}'$ . It depends on the parameter vector  $\boldsymbol{\theta}$ , which needs to be estimated.
- The spatial variation in  $w(\mathbf{s})$  is determined by the spatial correlation function  $\rho(\mathbf{s}, \mathbf{s}', \boldsymbol{\theta})$ . Locations that are closer to each other will tend to have more similar values of  $w(\mathbf{s})$  due to the positive correlation described by the spatial correlation function.
- It's important to note that, conditionally on  $w(\mathbf{s})$ , the observed data  $y_i(\mathbf{s}_i)$  are assumed to be mutually independent random variables. This means that the correlation in  $w(\mathbf{s})$  does not directly affect the correlation in the residuals  $\varepsilon(\mathbf{s})$ .
- $w(\mathbf{s})$  is spatially varying (in general each location  $\mathbf{s}$  has a different  $w(\mathbf{s})$  value)
- **How overfitting is avoided?** we can't use  $w(\mathbf{s})$  in order to estimate the exact value of  $y(\mathbf{s})$

$$y(\mathbf{s}) = \mathbf{x}(\mathbf{s})'\boldsymbol{\beta} + \alpha w(\mathbf{s}) + \varepsilon(\mathbf{s}) \quad \Rightarrow \quad \text{overfitting}$$

- The correlation function  $\rho(\mathbf{s}, \mathbf{s}', \boldsymbol{\theta})$  can be any positive-definite function. This condition imposes that the linear combination:

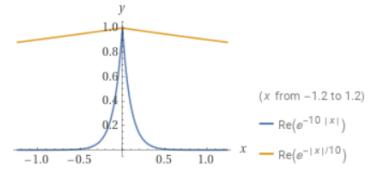
$$\sum_{i=1}^m \alpha_i \cdot w(\mathbf{s}_i) \quad \text{has positive variance}$$

- In most cases  $\rho(\mathbf{s}, \mathbf{s}', \boldsymbol{\theta}) = \rho(\|\mathbf{s} - \mathbf{s}'\|; \boldsymbol{\theta}) = \rho(u; \boldsymbol{\theta})$ , where  $\|\cdot\|$  is the distance (Euclidean or geodetic) between  $\mathbf{s}$  and  $\mathbf{s}'$
- In this case the spatial correlation only **depends on the distance**  $u$  and not on the coordinates.
- Instead if the correlation **depends on the coordinates**:  $w(0)$  and  $w(1)$  would be different from the correlation between  $w(1)$  and  $w(2)$

## 1.4 Exponential spatial correlation function

$$\rho(u; \theta) = \exp\left(-\frac{|\mathbf{s} - \mathbf{s}'|}{\theta}\right) = \exp\left(-\frac{u}{\theta}\right)$$

$\theta$  fixed, depends on distance  $u$



Observations:

- This function describes how the spatial correlation between two locations decreases exponentially as the distance between them, denoted as  $u$  increases.
- $\theta > 0$  is a scalar value. The larger the value of  $\theta$ , the slower the correlation decreases with distance, and the larger the effective range of spatial dependence. Conversely, a smaller value of  $\theta$  results in a faster decay of correlation with distance and a shorter effective range.
- The exponential correlation function is **simple** and computationally tractable, making it a common choice in spatial modeling. However, its simplicity can also be a limitation because it may **not be flexible enough** to capture complex spatial correlation patterns that are observed in some datasets.
- The suitability of the exponential correlation function depends on the specific characteristics of the spatial data under consideration. It is particularly appropriate when there is a clear spatial decay in correlation with distance, and the range of spatial dependence is relatively short. In cases where the spatial correlation pattern is not well described by exponential decay, alternative correlation functions (e.g., Gaussian, Matérn, or spherical) may be more appropriate.

## 1.5 Matérn spatial correlation function

$$\rho(u; \theta) = \{2^{K-1}\Gamma(K)\}^{-1} \left(\frac{u}{\phi}\right) K_K \left(\frac{u}{\phi}\right)$$

Observations:

- $K_K$ : is the modified Bessel function of the second kind of order  $K$
- $K$ : The parameter  $K$  is known as the smoothness or differentiability parameter. It determines how smooth or rough the spatial correlation function is. It is a positive real number  $K > 0$ .
- $u$ : This is the separation distance between two spatial points for which you want to compute the correlation. It is a non-negative real number.
- $\phi$ : with  $\phi > 0$  is the range parameter. It controls how quickly the spatial correlation decreases with increasing separation distance. A smaller  $\phi$  leads to a more rapid decrease in correlation with distance.

## 1.6 Gaussian process (GP)

- $w(\mathbf{s})$  not directly observed but can be inferred from the observed data set  $\{(\mathbf{x}_1(\mathbf{s}_1), y_1(\mathbf{s}_1)), \dots, \}$
- First, the distribution of  $w(\mathbf{s})$  must be specified in order to estimate it from the dataset
- Formally we want that  $w(\mathbf{s})$  is a zero mean Gaussian process (GP) with unitary variance and spatial correlation function  $\rho(u; \theta)$
- $w(\mathbf{s})$  is a gaussian process if the joint distribution of  $w(\mathbf{s}_1), \dots, w(\mathbf{s}_n)$  is multivariate normal (see SMS2) for any integer  $n$  and any set of locations  $(\mathbf{s}_1, \dots, \mathbf{s}_n)$
- 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 refers to the centres of the pixels

## 1.7 Matlab exercises 2-3

### 1.8 Spatial correlation - Regular grid

```

1 % Vettore coordinate
2 x = 0:1:99;
3 y = 0:1:99;
4
5 % Create a grid of coordinates (X, Y) using 'x' and 'y'vectors.
6 [X,Y] = meshgrid(x,y);
7
8 % Convert the grid coordinates (X, Y) into a 2D array 'coord'.
9 coord = [X(:) Y(:)];
10
11 % Plot the coordinates as points.
12 % Set the aspect ratio of the plot to be equal.
13 plot(coord(:,1),coord(:,2),'.');
14 axis equal
15
16 % Set a correlation length parameter.
17 theta = 0.1;
18
19 % Calculate pairwise Euclidean distances between points.
20 % distance matrix
21 dist = pdist2(coord,coord,"euclidean");

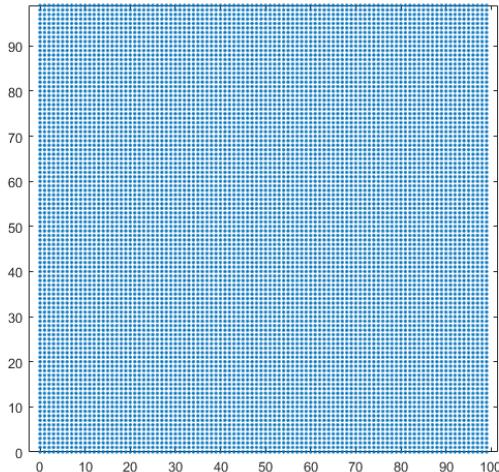
```

```

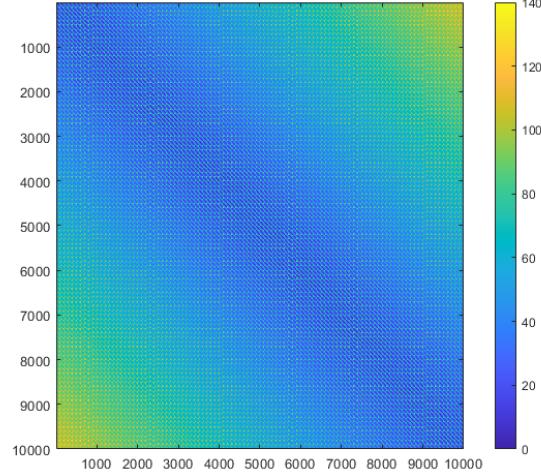
1 % Display the distance matrix as an image.
2 imagesc(dist);
3
4 % Calculate the spatial correlation using the exponential
5 % correlation function.
6 sp_corr = exp(-dist/theta);
7
8 % Display the spatial correlation matrix as an image.
9 imagesc(sp_corr);
10
11 % GP simulation:
12 % Generate a sample from a multivariate Gaussian distribution.
13 v = mvnrnd(zeros(1,size(dist,1)),sp_corr,1);
14
15 % Reshape the sample to match the grid shape.
16 v = reshape(v,size(X));
17
18 imagesc(v) % Display the generated GP sample as an image.
19 clim([-3,3]) % Set the color limit for the plot.
20 colorbar % Add a color bar to the plot.

```

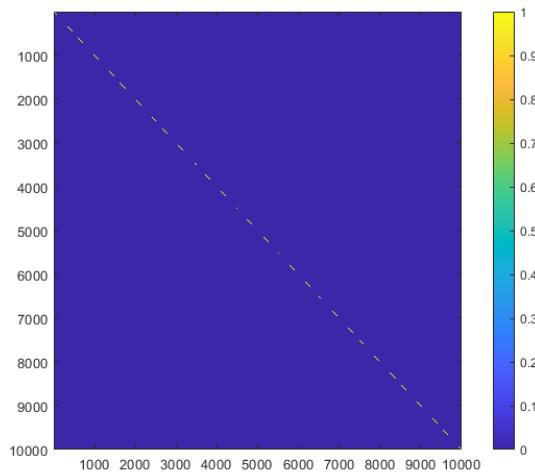
Uniform grid of the vectors  $\mathbf{X}$  and  $\mathbf{Y}$



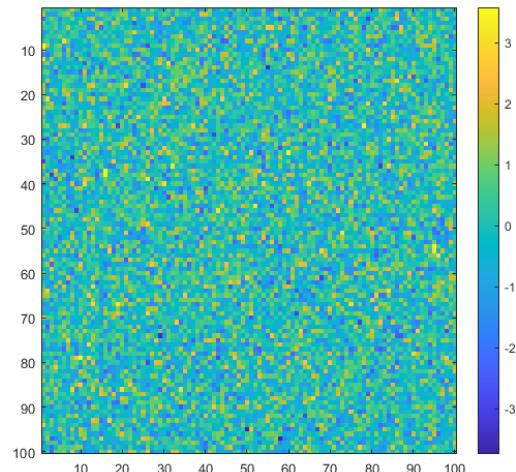
Distances matrix



Spatial correlation with  $\theta = 0.1$



GP sample as an image



## 1.9 Spatial correlation - Irregular grid

```

1 % Generate 200 random points in a 2D space with uniform coordinates
% between 0 and 1.
2 coord = unifrnd(0,1,200,2);
3
4 % Plot the irregularly spaced points.
5 plot(coord(:,1),coord(:,2),'.');
6
7 % Set the aspect ratio of the plot to be equal.
8 axis equal
9
10 % Set a correlation length parameter.
11 theta = 0.5;
12

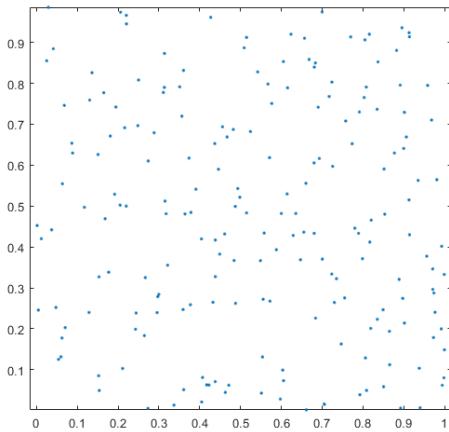
```

```

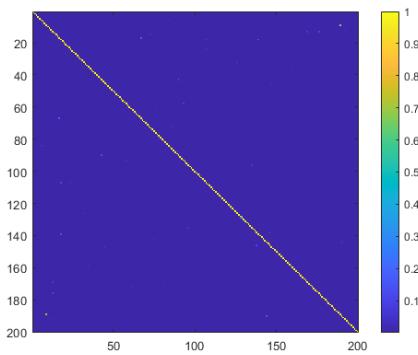
1 % Calculate pairwise Euclidean distances between points.
2 dist = pdist2(coord,coord,"euclidean");
3
4 % Display the distance matrix as an image.
5 imagesc(dist);
6
7 % Calculate the spatial correlation using the exponential
% correlation function
8 sp_corr = exp(-dist/theta);
9
10 % GP simulation:
11 % Generate a sample from a multivariate Gaussian distribution.
12 v = mvnrnd(zeros(1,size(dist,1)),sp_corr,1);
13
14 % Scatter plot with color-coded values based on the GP sample.
15 scatter(coord(:,1),coord(:,2),40,v,"filled")
16 clim([-3,3])
17 colorbar

```

Irregular grid of the vectors  $\mathbf{X}$  and  $\mathbf{Y}$

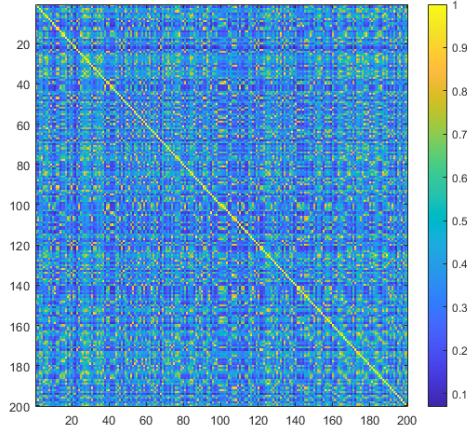


Spatial correlation with  $\theta = 0.006$

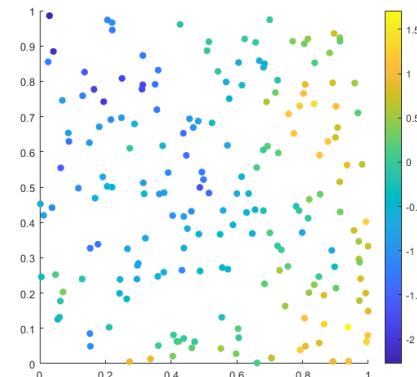
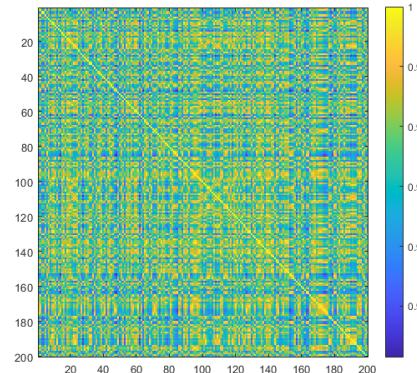


GP sample as an image

Distances matrix



Spatial correlation with  $\theta = 10.5$



## 2 Lesson 2

### 2.1 Spatial model with latent variable

$$y(\mathbf{s}) = \mathbf{x}(\mathbf{s})'\boldsymbol{\beta} + \alpha w(\mathbf{s}) + \varepsilon(\mathbf{s}) \quad (1)$$

Remarks:

- $w(\mathbf{s}) \sim GP(0, \rho(||\mathbf{s} - \mathbf{s}'||; \boldsymbol{\theta}))$ , is an approximation, in IRL we don't really know its distribution. Assuming that means if there is etheroschedasticity the variance not captured is inside the  $\varepsilon$
- $w(s)$  is a gaussian process and its not a normal multivariate in fact the mean is a scalar and the **variance** is not a matrix but is a **function** related to the correlation function and parameterized by  $\boldsymbol{\theta}$
- where  $||\mathbf{s} - \mathbf{s}'||$  is the euclidean distance between 2 generic points
- $\rho(||\mathbf{s} - \mathbf{s}'||; \boldsymbol{\theta}) = corr(w(\mathbf{s}), w(\mathbf{s}'))$
- $\varepsilon(\mathbf{s}) \sim N(0, \sigma_\varepsilon^2)$ , we can use the normal because there is no correlation in space and the variables are independent
- The unknown parameter set is  $\Psi = \{\boldsymbol{\beta}, \alpha, \sigma_\varepsilon^2, \boldsymbol{\theta}\}$
- How to estimate  $\Psi$  from data? **MLE** → **EM**

**Remarks:**  $w(\mathbf{s}) \sim GP$  is an approximation, we dont really know what is its distribution. In fact if we impose that  $w(s) \sim GP$  and  $w(s)$  is not a gaussian process,  $\mathbf{x}(\mathbf{s})'\boldsymbol{\beta}$  are not capable to explain well the data observed and the **error** is in  $\varepsilon(\mathbf{s})$ . We can confirm this if:

- there si spatial correlation left inside  $\varepsilon(\mathbf{s})$ .
- or  $\varepsilon(\mathbf{s})$  is etheroschedastic in space

### 2.2 Maximum likelihood estimate

- We rely on MLE to estimate the model parameter vector  $\Psi$ , because MLE has good properties and in case we can use the EM algorithm (if needed).
  - **Correct Model Specification:** the model you choose should be a valid representation of the data-generating process. If the model is misspecified, MLE estimates may not be accurate.
  - **the observations must be IID:** assumes that the data points are drawn from the same probability distribution with the same parameter values and the observations are independent of each other. This means that the probability of one observation does not depend on the values of other observations.
  - **Finite Parameter Space:** MLE assumes that the parameter space is finite. This means that the parameter values being estimated exist within a specific, well-defined range.
  - We use **the Expectation-Maximization (EM)** algorithm when the data has missing or latent variables.
- with EM the missing values can be:
  - $y(\mathbf{s})$  but if we have a spatial dataset we just delete that information, instead in spatial temporal dataset we dont delete the  $y(t)$  but assign it a NA value
  - or  $w(s)$
- Likelihood function of (3.2) is

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

- $\mathbf{y} = (y_1, \dots, y_n)'$  is the vector of observations at  $n$  spatial locations
- $\mathbf{w} = (w_1, \dots, w_n)'$  is the vector of latent variables at  $n$  spatial locations ⇒ **we pretend for the moment to observe it**
- $\mathbf{X}$  is the  $n \times p$  design matrix, with  $p$  covariates
- we **assume** that the covariates does not influence the  $w(s)$ , the covariates are numbers without uncertainty

## 2.3 Likelihood decomposition

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

Remarks:

- $L(\Psi; \mathbf{y}, \mathbf{w}, \mathbf{X})$  is the complete-data likelihood (which assumes  $\mathbf{w}$  to be known)
- In reality I don't observe  $\mathbf{w}$  so I'm forced to do this decomposition
- $L(\beta, \alpha, \sigma_\varepsilon^2; \mathbf{y} | \mathbf{x}, \mathbf{X})$  and  $L(\boldsymbol{\theta}; \mathbf{w})$  are described by a normal n-variate distributions. This derive from the fact that:

$$w(s) \sim GP \text{ and } \varepsilon(s) \sim N$$

## 2.4 Log-likelihood function

$$-2\log(L_\Psi) = \log|\Sigma_\varepsilon| + \mathbf{e}' \Sigma_\varepsilon \mathbf{e} + \log|\Sigma_w| + \mathbf{w}' \Sigma_w^{-1} \mathbf{w}$$

Remarks:

- As usual we prefer to work with  $\log(L_\Psi) \rightarrow$  maximize
- $-2\log(L_\Psi) \rightarrow$  minimize
- $\mathbf{e} = \mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \alpha\mathbf{w}$ , error vector (we still pretend to know  $\mathbf{w}$ )
- $\Sigma_\varepsilon = \sigma_\varepsilon^2 \mathbf{I}_n$ , with  $\mathbf{I}_n$  the identity matrix of dimension  $n \rightarrow$  diagonal matrix
- $\Sigma_\varepsilon$  we don't have spatial correlation that means  $\varepsilon$  is independent so we only have a diagonal matrix with the same value on the diagonal and all the covariances are 0 (omoschedasticity, same variance between each point in space  $\rightarrow \sigma_\varepsilon^2$ )
- $\Sigma_w$  is the  $n \times n$  correlation matrix (e.g.  $\exp(-D/\theta)$ , with  $D$  the distance matrix)
- If we know  $\mathbf{w}$  we can minimize directly Log-likelihood (only one unknown)

## 2.5 ML estimate

$$\hat{\Psi} = \underset{\boldsymbol{\beta}, \alpha, \sigma_\varepsilon^2, \boldsymbol{\theta}}{\operatorname{argmin}} \log|\Sigma_\varepsilon| + \mathbf{e}' \Sigma_\varepsilon \mathbf{e} + \log|\Sigma_w| + \mathbf{w}' \Sigma_w^{-1} \mathbf{w}$$

Remarks:

- Argmin because we are considering  $-2\log(L_\Psi)$
- Unfortunately minimizing  $-2\log(L_\Psi)$  is not feasible (not computationally applicable), plus  $\mathbf{w}$  is **latent and not observed (we can't calculate this parts)**  $\Rightarrow$  we must rely on the EM algorithm

## 2.6 EM algorithm

- The EM is an iterative algorithm for MLE and under some assumptions the estimate of EM is the same as the ML estimate  $\Rightarrow$  this is the reason why we use it
- First iteration starts with fixed initial values chosen by us  $\hat{\Psi}^0$  (usually given by OLS and method of moments). If we initialize  $\hat{\Psi}^0$  with parameters that are really far from the optimal there is the risk to converge to **local minimum**
- EM algorithm has two steps: Expectation and Maximization

### E-step:

$$Q(\Psi, \hat{\Psi}^m) = \mathbb{E}_{\hat{\Psi}^m} [-2\log L(\Psi; \mathbf{y}, \mathbf{w}, \mathbf{X}|\mathbf{y})]$$

### M-step (Maximization):

Closed form:

$$\begin{aligned} Q(\Psi, \hat{\Psi}^m) &= \mathbb{E}_{\hat{\Psi}^m} [-2\log L(\Psi; \mathbf{y}, \mathbf{w}, \mathbf{X}|\mathbf{y})] \\ &= \text{tr}[\Sigma_{\varepsilon}^{-1} (\mathbb{E}(\mathbf{e}|\mathbf{y}) \mathbb{E}(\mathbf{e}|\mathbf{y})' + \text{Var}(\mathbf{e}|\mathbf{y}))] \\ &\quad + \text{tr}[\Sigma_{\varepsilon}^{-1} (\mathbb{E}(\mathbf{w}|\mathbf{y}) \mathbb{E}(\mathbf{w}|\mathbf{y})' + \text{Var}(\mathbf{w}|\mathbf{y}))] \end{aligned}$$

$$\hat{\Psi}^{(m+1)} = \arg\max_{\Psi} Q(\Psi, \hat{\Psi}^{(m)})$$

$$\frac{dQ(\Psi, \hat{\Psi}^{(m)})}{d\Psi} = 0$$

Remarks:

1.  $\mathbb{E}(\mathbf{e}|\mathbf{y}) = \mathbf{y} - \mathbf{X}\beta - \alpha\mathbb{E}(\mathbf{w}|\mathbf{y})$
2.  $\text{Var}(\mathbf{e}|\mathbf{y}) = \text{Var}(\mathbf{y} - \mathbf{X}\beta - \alpha\mathbf{w}|\mathbf{y}) = \alpha^2 \text{Var}(\mathbf{w}|\mathbf{y})$
3.  $\mathbb{E}(\mathbf{w}|\mathbf{y}) = \text{Cov}(\mathbf{w}, \mathbf{y}) \text{Var}(\mathbf{y})^{-1} [\mathbf{y} - \mathbf{X}\beta]$
4.  $\text{Var}(\mathbf{w}|\mathbf{y}) = \Sigma_w - \text{Cov}(\mathbf{w}, \mathbf{y}) \text{Var}(\mathbf{y})^{-1} \text{Cov}(\mathbf{w}, \mathbf{y})'$
5.  $\text{Var}(\mathbf{y}) = \text{Var}(\mathbf{X}\beta + \alpha\mathbf{w} + \boldsymbol{\varepsilon}) = \text{Var}(\alpha\mathbf{w} + \boldsymbol{\varepsilon})$   
 $= \alpha^2 \text{Var}(\mathbf{w}) + \text{Var}(\boldsymbol{\varepsilon}) + 2\text{Cov}(\mathbf{w}, \boldsymbol{\varepsilon})$
6.  $\text{Var}(\mathbf{w}) = \Sigma_w$
7.  $\text{Var}(\boldsymbol{\varepsilon}) = \Sigma_{\varepsilon} = \sigma_{\varepsilon}^2 \mathbf{I}_n$
8.  $2\text{Cov}(\mathbf{w}, \boldsymbol{\varepsilon}) = \mathbf{0}$
9.  $\text{Cov}(\mathbf{w}, \mathbf{y}) = \text{Cov}(\mathbf{w}, \mathbf{X}\beta + \alpha\mathbf{w} + \boldsymbol{\varepsilon}) = \text{Cov}(\mathbf{w}, \alpha\mathbf{w})$   
 $= \alpha\text{Cov}(\mathbf{w}, \mathbf{w}) = \alpha\text{Var}(\mathbf{w}) = \alpha\Sigma_w$

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

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

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

$$\boldsymbol{\theta}^{(m+1)} = \arg\min_{\boldsymbol{\theta}} \log |\Sigma_w^{-1}(\boldsymbol{\theta})| + \text{tr}[\Sigma_w^{-1}(\boldsymbol{\theta})(\widehat{\mathbf{w}}\widehat{\mathbf{w}}')]$$

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

### Importante:

- (1) Il valore atteso di un vettore di numeri  $\mathbf{y}$  sono i numeri stessi e il valore atteso si riduce su  $\mathbf{w}$ . Sto cercando di ottenere qual'è il valore atteso della variabile latente  $\mathbf{w}$  condizionatamente ai dati che ho osservato  $\mathbf{y}$ .
- (2) La varianza condizionata di un scalare è 0 e sopravvive la var su  $\mathbf{w}$ . Fare l'inversa di una matrice è computazionalemente oneroso  $\text{Var}(\mathbf{y})^{-1}$
- (3) rappresenta la stima della variabile latente. Se la  $\text{Var}(\mathbf{y}) \rightarrow \infty$ , ossia l'errore di misura è molto alto  $\text{Var}(\boldsymbol{\varepsilon})$  non posso fidarmi di  $[\mathbf{y} - \mathbf{X}\beta]$  (errore di modello) perché è tutto errore e dunque la stima di  $\mathbf{w}$  è 0, al contrario se  $\text{Var}(\mathbf{y}) \rightarrow \text{Var}(\mathbf{w})$ , dunque non abbiamo errore di misura e quindi nemmeno errore di modello e la stima di  $\mathbb{E}(\mathbf{w}|\mathbf{y}) = \mathbf{y} - \mathbf{X}\beta$

$$\begin{array}{lll} \text{Var}(\boldsymbol{\varepsilon}) = \infty & \text{Var}(\mathbf{y}) \rightarrow \infty & \mathbb{E}(\mathbf{w}|\mathbf{y}) = 0 = \text{mean GP} \\ \text{Var}(\boldsymbol{\varepsilon}) = 0 & \text{Var}(\mathbf{y}) \rightarrow \text{Var}(\mathbf{w}) & \mathbb{E}(\mathbf{w}|\mathbf{y}) = \mathbf{y} - \mathbf{X}\beta \end{array}$$

- (4) viene interpretata come incertezza

$$\begin{array}{lll} \text{Var}(\boldsymbol{\varepsilon}) = \infty & \text{Var}(\mathbf{y}) \rightarrow \infty & \text{Var}(\mathbf{w}|\mathbf{y}) = \Sigma_w \\ \text{Var}(\boldsymbol{\varepsilon}) = 0 & \text{Var}(\mathbf{y}) \rightarrow \text{Var}(\mathbf{w}) & \text{Var}(\mathbf{w}|\mathbf{y}) = \Sigma_w - \Sigma_w = 0 \end{array}$$

- (7) Matrice delle distanze è deterministica,  $\boldsymbol{\theta}$  è un parametro
- (8) Assunzione del modello, le variabili tra di loro sono indipendenti
- (9)  $\text{Cov}(\mathbf{w}, \mathbf{X}\beta) = 0$  perché  $\mathbf{X}$  costante,  $\text{Cov}(\mathbf{w}, \boldsymbol{\varepsilon}) = 0$ , rimane  $\text{Cov}(\mathbf{w}, \alpha\mathbf{w}) = 0$
- ( $\alpha^{m+1}$ ) dipende dalla (3) è possano esserci problemi di identificabilità, var alta  $\alpha$  non ha significato
- ( $\boldsymbol{\beta}^{m+1}$ ) aggiornamento di beta, togliamo alpha appena stimato e il valore atteso di  $\mathbf{w}$ , perché non voglio avere l'effetto di  $\mathbf{w}$  su  $\boldsymbol{\beta}$

- $(\sigma_{\varepsilon}^{m+1})$  il residuo ora é  $\mathbf{e}|\mathbf{y} = \mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \alpha\mathbb{E}(\mathbf{w}|\mathbf{y})$
- $(\boldsymbol{\theta}^{m+1})$  é l'unica quantitá che non riesco a stimare in forma chiusa, infatti questa stima é data da una minimizzazione della verosimiglianza che riguarda solo theta. Il minimizzatore cambia theta internamente e ogni volta che lo fa cambiano anche le matrici coinvolte  $\Rightarrow$  **penalizzazione computazionale molto elevata**
- Se ho dei dubbi che la mia stima del EM non sia quella della ML, posso randomizzare i valori iniziali del mio  $\hat{\Psi}^0$  e vedo se l'algoritmo converge alle stesse soluzioni
- Algoritmo EM non fornisce direttamente alcuna incertezza sui parametri stimati:  $\alpha, \boldsymbol{\beta}, ..$  Mi da solo l'incertezza della variabile lantente  $Var(\mathbf{w}|\mathbf{y})$
- Solitamente l'incertezza sui parametri tende ad essere molto bassa, tranne sul  $(\boldsymbol{\theta}^{m+1})$ , questo perché risulta essere poco identificabile, anche raddoppiandolo la verosimiglianza cambia poco
- $\mathbf{w}$  viene compensato da  $\mathbf{X}\boldsymbol{\beta}$ , quindi anche cambiandolo é il modello che mi da una mano
- Come da variogramma:  $\theta$  piú é grande piú c'è correlazione spaziale, se dunque mando  $\theta \rightarrow \infty$  ho  $\mathbf{w}$  costante, ossia la correlazione spaziale é cosí forte che indipendentemente dalla distanza  $\mathbf{w}$  rimane costante

## 2.7 Variogram

- In spatial statistics, a variogram (also known as a semivariogram) is a graphical and mathematical tool used to quantify spatial dependence or spatial autocorrelation in a set of spatial data
- Spatial dependence refers to the idea that the values of a variable at one location are related to the values of the same variable at nearby locations. The variogram helps to characterize the spatial structure and variability of a spatial process.
- The variogram is typically calculated using pairs of observations at different locations within a study area. **The variogram measures the variability between pairs of points as a function of the distance or lag separating them.**
- The variogram is denoted as  $2\gamma(\mathbf{s}, \mathbf{s}')$ , instead the semivariogram is half the variogram  $\gamma(\mathbf{s}, \mathbf{s}')$ .

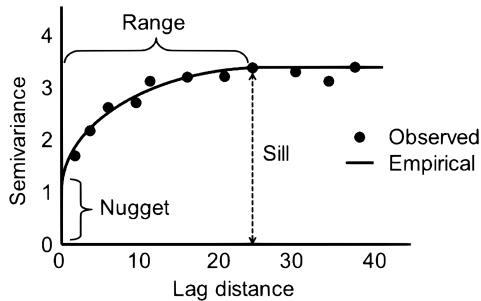
The semivariogram was defined as half the average squared difference between the values at points ( $\mathbf{s}_1$  and  $\mathbf{s}_2$ ) separated at distance  $\mathbf{h}$ , where  $\mathbf{M}$  is a point in the geometric field  $\mathbf{V}$  and  $f(\mathbf{M})$  is the value at that point. The triple integral is over 3 dimensions, instead  $\mathbf{h}$  is the separation distance (meters or km).

*Example: the value  $f(M)$  could represent the iron content in soil at some location  $M$  (with geographic coordinates of latitude, longitude, and elevation) over some region  $V$  with element of volume  $dV$ . To obtain the semivariogram for a given  $\gamma(h)$ , all pairs of points at that exact distance would be sampled. In practice it is impossible to sample everywhere, so the empirical variogram is used instead.*

The variogram is twice the semivariogram and can be defined, equivalently, as the **variance** of the difference between field values at two locations ( $\mathbf{s}_1$  and  $\mathbf{s}_2$ , note change of notation from  $M$  to  $\mathbf{s}$  and  $f$  to  $Z$ ) across realizations of the field

$$\begin{aligned} \text{semivariogram:} \quad \gamma(h) &= \frac{1}{2V} \int \int \int_V [f(\mathbf{M} + \mathbf{h}) - f(\mathbf{M})]^2 dV \\ \text{variogram:} \quad 2\gamma(h) &= var \left( (Z(\mathbf{s}_1) - \mu) - (Z(\mathbf{s}_2) - \mu) \right)^2 \end{aligned}$$

- The variogram is then plotted as a function of lag distance, showing how the semivariance changes with increasing separation between points. The resulting variogram plot can provide insights into the spatial structure of the variable under consideration, helping analysts to understand patterns such as clustering or spatial trends.



### Parameters:

- *nugget n*: represents the discontinuity or the "jump" of the semivariogram at the origin (lag distance equals zero). It accounts for the variability at very short distances that might be due to measurement errors or other factors not captured by the spatial model. A higher nugget indicates a greater level of variability at very short distances.
- *sill s*: The sill is the limiting value of the variogram as the lag distance tends to infinity. It represents the total variability or the maximum semivariance in the spatial process. A higher sill indicates a larger overall variability in the spatial process.
- *range r*: The distance in which the difference of the variogram from the sill becomes negligible. It signifies the distance at which spatial autocorrelation or dependence between data points is no longer significant. A longer range indicates that spatial correlation or similarity between points persists over larger distances.

**Nugget**      *Represents short-range variability or variability at very small distances*

**Sill**            *Represents the total variability or the maximum semivariance in the spatial process*

**Range**          *Represents the distance at which spatial correlation becomes negligible*

### Empirical variogram

- Variogram is made on countinous data which in IRL we don't have so the empirical variogram is base on the given data we have
- Calculate the differences between all pairs of data points at different lag distances (distances between points) and for each lag distance, calculate the semivariance, which indicates how much values differ at a given distance
- Group the semivariance values into bins based on the lag distance and calculate the average semivariance within each bin
- Plot the lag distance on the x-axis and the average semivariance on the y-axis. The resulting plot is the empirical variogram, which provides insights into the spatial structure of the data.

### Applications:

- It a model of the temporal/spatial correlation of the observed phenomenon
- The **experimental variogram** that is a visualisation of a possible spatial/temporal correlation and the **variogram model** that is further used to define the weights of the kriging function
- The experimental variogram is an empirical estimate of the covariance of a Gaussian process. As such, it may not be positive definite and hence not directly usable in kriging, without constraints or further processing.

### 3 Lesson 3

#### 3.1 Spatial model with latent variable

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

Remarks:

- $w(\mathbf{s}) \sim GP(0, \rho(||\mathbf{s} - \mathbf{s}'||; \boldsymbol{\theta}))$
- $\rho(||\mathbf{s} - \mathbf{s}'||; \boldsymbol{\theta}) = corr(w(\mathbf{s}), w(\mathbf{s}'))$
- $\varepsilon(\mathbf{s}) \sim N(0, \sigma_\varepsilon^2)$  measurement/model error
- $\mathbf{x}(\mathbf{s})'\boldsymbol{\beta}$  are the regressors associated with their coefficients.  $\mathbf{x}(\mathbf{s})$  are numbers with no variance and no distribution, instead  $\boldsymbol{\beta}$  is the related coefficients estimated with EM.
- The unknown parameters  $\Psi = \{\boldsymbol{\beta}, \alpha, \sigma_\varepsilon^2, \boldsymbol{\theta}\}$  is estimated using the EM algorithm

#### 3.2 Prediction for one site $s_i$

$$\hat{y}(\mathbf{s}_i) = \mathbf{x}(\mathbf{s}_i)' \hat{\boldsymbol{\beta}} + \hat{\alpha} \hat{w}(\mathbf{s}_i)$$

Remarks:

- The above formula gives the prediction at spatial locations  $s_i$ , with  $i = 1, \dots, n$
- Seems to be more like a **filtering** than a prediction. Measurement error is cutted off
- We are making the prediction in a specific space location  $s_i$
- Monitoring stations in space → we want to estimate the pollutants concentration → regular grid.  $s_i$  will be a pixel of the grid, the center of the pixel
- $\hat{w}(s_i) = \mathbb{E}[w(s_i)|Y]$ , we estimate  $\hat{w}$  using the observations in the dataset  $Y$
- $\hat{y}(s_i)$  the prediction in a specific location requires the regressors in that location  $\mathbf{x}(\mathbf{s}_i)' \hat{\boldsymbol{\beta}}$  if not we can't do the prediction. A way to avoid this is to interpolate across space and use the interpolated values as  $\mathbf{x}(\mathbf{s}_i)$  in the model. This is not the best approach because  $\mathbf{x}(\mathbf{s}_i)$  should be a value with no variance, so deterministic, the interpolation will add an error to the model.

Important:

- First of all we need  $\mathbf{x}(\mathbf{s})$  (at spatial covariates  $\mathbf{s}$ )
- $\mathbb{E}(\mathbf{w}|\mathbf{y}) = Cov(\mathbf{w}, \mathbf{y})Var(\mathbf{y})^{-1}[\mathbf{y} - \mathbf{X}\boldsymbol{\beta}]$
- Similarly  $\hat{w}(\mathbf{s}) = \mathbb{E}(w(\mathbf{s})|\mathbf{y}) = Cov(w(\mathbf{s}), \mathbf{y})Var(\mathbf{y})^{-1}[\mathbf{y} - \mathbf{X}\boldsymbol{\beta}] \in \mathbb{R}$
- In practice the prediction  $\hat{w}(\mathbf{s})$  depends on the vector  $\mathbf{y}$  of all observations
- More in details:

$$\begin{aligned} \hat{w}(\mathbf{s}) &= \mathbb{E}(w(\mathbf{s})|\mathbf{y}) = Cov(w(\mathbf{s}), \mathbf{y})Var(\mathbf{y})^{-1}[\mathbf{y} - \mathbf{X}\boldsymbol{\beta}] \\ Cov(w(\mathbf{s}), \mathbf{y}) &= Cov(w(\mathbf{s})|\mathbf{X}\boldsymbol{\beta} + \alpha\mathbf{w} + \varepsilon) \\ &= Cov(w(\mathbf{s}), \alpha\mathbf{w}) \\ &= \alpha Cov(w(\mathbf{s}), \mathbf{w}) \in \mathbb{R}^{1 \times n} \\ Cov(w(\mathbf{s}), w(\mathbf{s}_i)) &= exp\left(-\frac{||\mathbf{s} - \mathbf{s}_i||}{\theta}\right) \end{aligned}$$

- $Var(\mathbf{y})^{-1}[\mathbf{y} - \mathbf{X}\boldsymbol{\beta}]$  seen in the EM algorithm and do not depend on  $w(\mathbf{s})$
- $Cov(\mathbf{w}(\mathbf{s}), \mathbf{y})$  can be seen as a weight vector, so very far from the network this covariance calculated by the exponential will be close to 0 and the prediction we are making will be as well close to 0

### 3.3 Spatial prediction for multiple sites $S$

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

- For one site  $\mathbf{s}_i$ , we try to cover the region with the grid and we predict the pollutant concentration for each pixel of the grid, instead the prediction for multiple sites, we collect all the locations into  $\mathbf{S}$  and the prediction will be a vector. This time we are doing a multiple prediction on multiple points on the grid
- Also if we want we can do predictions on blocks of the grid and not on all the grid
- Prediction is done for multiple spatial locations  $\mathbf{S} = \{\mathbf{s}_1, \dots, \mathbf{s}_M\}$
- $Cov(\mathbf{w}(\mathbf{S}), \alpha\mathbf{w}) = \alpha Cov(\mathbf{w}(\mathbf{S}), \mathbf{w})$
- $\alpha Cov(\mathbf{w}(\mathbf{S}), \mathbf{w})$  is a  $M \times n$  matrix  $\rightarrow$  computation problems when  $M$  is big
- $Var(\mathbf{y})^{-1}[\mathbf{y} - \mathbf{X}\boldsymbol{\beta}]$  is computed only one time, for one site we should have computed this operation for each point on the grid
- $Var(\mathbf{w}(\mathbf{S})|\mathbf{y})$  is the prediction uncertainty (only the diagonal, outside we have the covariances)

### 3.4 Multivariate models

$$\mathbf{y}(\mathbf{s}) = \mathbf{X}(\mathbf{s})'\boldsymbol{\beta} + \boldsymbol{\alpha} \odot \mathbf{w}(\mathbf{s}) + \boldsymbol{\varepsilon}(\mathbf{s})$$

$$\mathbf{y}(\mathbf{s}) = \mathbf{X}(\mathbf{s})'\boldsymbol{\beta} + \mathbf{A}\mathbf{w}(\mathbf{s}) + \boldsymbol{\varepsilon}(\mathbf{s})$$

Remarks:

- $\mathbf{y}(\mathbf{s})$  is a  $p \times 1$  vector, where  $p$  are the number of variables we want to model
- $\mathbf{w}(\mathbf{s}) \sim GP_p(\mathbf{0}, \mathbf{V}\rho(||\mathbf{s} - \mathbf{s}'||; \boldsymbol{\theta}))$  is a p-variate Gaussian random process
- $\mathbf{V}$  is a  $p \times p$  correlation matrix
- $\boldsymbol{\varepsilon}(\mathbf{s}) \sim N_p(\mathbf{0}, \boldsymbol{\Sigma}_{\varepsilon}^2)$  is a p-variate Normal random variable with  $\boldsymbol{\Sigma}_{\varepsilon}^2$  diagonal
- The unknown parameter set is  $\Psi = \{\boldsymbol{\beta}, \boldsymbol{\alpha}, \boldsymbol{\Sigma}_{\varepsilon}^2, \boldsymbol{\theta}, \mathbf{V}\}$

Importante:

- Voglio modellare congiuntamente più variabili di output nello stesso modello, ossia più  $\mathbf{y}$
- Stimiamo un modello unico che modellizza tutte le  $p$  variabili e non  $p$  modelli univariati
- Conviene scegliere questo modello quando una delle variabili è osservata su pochi siti ma è correlata con le altre, in questo modo sfrutto la correlazione fra le variabili per aggiungere informazione al modello. Nel caso di incorrelazione non mi conviene fare il modello multivariato
- Le nostre covariate nel modello sono dei numeri deterministici senza incertezza, nel caso di incertezza non posso usarla come regressore. Inoltre il regressore dev'essere sempre presente per poter farlo usare come covariata del modello  $\rightarrow$  al massimo non uso alcuna covariata
- $\mathbf{X}$  è una matrice a blocchi non più un vettore, e ogni blocco si riferisce a ciascuna variabile del modello, ossia ciascuna variabile  $y$  ha il suo set di covariate che possono differire dalle altre
- $\boldsymbol{\beta}$  anch'esso sarà strutturato in sotto-vettori e ciascun sottovettore si riferisce ad una variabile  $y_i$

$p = 2$ , covariate di  $y_1$  e  $y_2$

$$\begin{bmatrix} \textcolor{red}{y_1} \\ \textcolor{blue}{y_2} \end{bmatrix} = \begin{bmatrix} (\textcolor{red}{x_{11}} & \textcolor{red}{x_{12}}) & 0 & 0 & 0 \\ 0 & 0 & (\textcolor{blue}{x_{21}} & \textcolor{blue}{x_{22}} & \textcolor{blue}{x_{23}}) \end{bmatrix} \begin{bmatrix} \textcolor{red}{\beta_1} \\ \textcolor{red}{\beta_2} \\ \textcolor{blue}{\beta_3} \\ \textcolor{blue}{\beta_4} \\ \textcolor{blue}{\beta_5} \end{bmatrix}$$

- $\mathbf{w}(\mathbf{s}) \sim GP_p(\mathbf{0}, \mathbf{V}_p \rho(||\mathbf{s} - \mathbf{s}'||; \boldsymbol{\theta}))$ : date  $p$  variabili esiste un  $GP$  per ciascuna variabile e probabilmente diversi tra di loro.
- $\mathbf{V}_p \cdot \rho$  matrice dove ogni elemento è una funzione e  $\mathbf{V}$  è una matrice di correlazione
- $\boldsymbol{\theta}$  parametro che mi gestisce la correlazione spaziale, le  $p$  variabili sono obbligate a condividere lo stesso  $\boldsymbol{\theta}$  per la funzione di correlazione  $\rho$ , questo può essere un problema perché le variabili posso essere correlate in modo diverso. Tale scelta viene fatta in modo tale che la matrice sia semi-definita positiva  $\mathbf{V}_p \cdot \rho$
- $\boldsymbol{\varepsilon}(\mathbf{s})$ : date  $p$  variabili su ciascuna variabile avrà un errore di misura differente e con scale differenti, inoltre la varianza risulta essere  $\boldsymbol{\Sigma}_\varepsilon^2$  diagonale, perché non ho covarianza tra gli errori, questo perché la correlazione la gestisce  $\mathbf{w} \rightarrow$  gli errori di misura sulle  $p$  variabili sono indipendenti
- $A$  è una matrice diagonale  $p \times p$

### 3.5 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)

*data observed n (sites) = 1000 → 1000 × 1000 matrices*

*with p = 5 → 5000 × 5000 matrices*

#### Variables observed on n sites

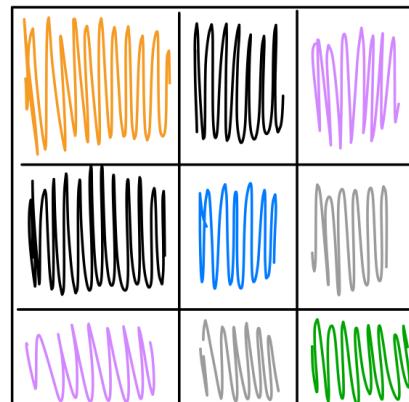
$$n_1 = 1000 \quad n_2 = 300 \quad n_3 = 150 \quad \rightarrow \quad n = 1450$$

*matrice di correlazione del  $\mathbf{w}$  dato da un  $GP_3$*

$\exp(-\frac{D_1}{\theta})$ $D_1 = 1000 \times 1000$	$\downarrow \exp(-\frac{D_{12}}{\theta})$ //	
//	$\exp(-\frac{D_{23}}{\theta})$ $D_2 = 300 \times 300$	
		$\exp(-\frac{D_3}{\theta})$ $D_3 = 150 \times 150$

Simmetrica  
ottica simmetrica

W



### 3.6 Bivariate model

- A bivariate model is a (simple) special case of the multivariate model.
- $\mathbf{x}_1(\mathbf{s})'$  and  $\mathbf{x}_2(\mathbf{s})'$  can have different lengths

$$\begin{bmatrix} y_1(\mathbf{s}) \\ y_2(\mathbf{s}) \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1(\mathbf{s})' & \mathbf{0} \\ \mathbf{0} & \mathbf{x}_2(\mathbf{s})' \end{bmatrix} \begin{bmatrix} \boldsymbol{\beta}_1(\mathbf{s}) \\ \boldsymbol{\beta}_2(\mathbf{s}) \end{bmatrix} + \begin{bmatrix} \alpha_1 w_1(\mathbf{s}) \\ \alpha_2 w_2(\mathbf{s}) \end{bmatrix} + \begin{bmatrix} \varepsilon_1(\mathbf{s}) \\ \varepsilon_2(\mathbf{s}) \end{bmatrix}$$

$$\mathbf{V} = \begin{bmatrix} 1 & \text{corr}(w_1(\mathbf{s}), w_2(\mathbf{s})) \\ \text{corr}(w_2(\mathbf{s}), w_1(\mathbf{s})) & 1 \end{bmatrix}$$

$$\boldsymbol{\Sigma}_{\varepsilon}^2 = \begin{bmatrix} \sigma_{1,\varepsilon}^2 & 0 \\ 0 & \sigma_{2,\varepsilon}^2 \end{bmatrix}$$

### 3.7 Linear model of coregionalization

- $\mathbf{V}\rho(||\mathbf{s} - \mathbf{s}'||; \boldsymbol{\theta})$  is called linear coregionalization model
- $\boldsymbol{\theta}$  is common to the  $p$  variables is really difficult to identify
- This could be a limit because all the variables are forced to share the same spatial correlation (same function and strength)
- $\mathbf{V}$  is a symmetric correlation matrix (diagonal elements equal to 1), only  $\frac{p(p-1)}{2}$  elements of  $\mathbf{V}$  are estimated

### 3.8 Data structure

A multivariate data set can be classified depending on spatial locations:

- **Isotopic:** all  $p$  variables are observed at the same  $n$  spatial locations.
  - There is a single distance matrix, as all variables share the same set of spatial locations.
  - $\boldsymbol{\alpha} \odot \mathbf{w}(\mathbf{s})$  I can directly execute the alpha operation inside  $\mathbf{w}(\mathbf{s})$  and I will have not a correlation matrix but a covariance one. Estimated in closed form
- **Fully heterotopic:** the  $p$  variables do not share a single spatial location.
  - There are  $p$  distance matrices, each corresponding to a different variable.
  - $\boldsymbol{\alpha} \odot \mathbf{w}(\mathbf{s})$  can't be estimated in closed form because the  $\mathbf{w}$  matrix is rectangular
- **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.
  - There are  $p$  distance matrices, each corresponding to a different variable. However, the matrices may have different dimensions based on the subset of spatial locations where each variable is observed.
  - $\boldsymbol{\alpha} \odot \mathbf{w}(\mathbf{s})$  can't be estimated in closed form because the  $\mathbf{w}$  matrix is rectangular

## 4 Lesson 4

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

$$y(\mathbf{s}, t) = \mathbf{x}_\beta(\mathbf{s}, t)' \boldsymbol{\beta} + \mathbf{x}_z(\mathbf{s})' \mathbf{z}(t) + \alpha w(\mathbf{s}, t) + \varepsilon(\mathbf{s}, t)$$

$$\mathbf{z}(t) = \mathbf{G}\mathbf{z}(t-1) + \boldsymbol{\eta}(t)$$

Remarks:

- $\mathbf{x}_\beta(\mathbf{s}, t)$  and  $\mathbf{x}_z(\mathbf{s})$  are vectors of covariates, they can't contain missing values
- $w(\mathbf{s}, t) \sim GP(0, \rho(||\mathbf{s} - \mathbf{s}'||; \boldsymbol{\theta}))$  is correlated over space but IID over time
- $\mathbf{z}(t)$  is  $q \times 1$  dimensional with Markovian dynamics
- $\mathbf{G}$  is a stable  $q \times q \Rightarrow$  transition matrix, it can be full, diagonal
- $\boldsymbol{\eta}(t) \sim N(\mathbf{0}, \Sigma_\eta)$  is the innovation with  $\Sigma_\eta$  the variance-covariance matrix
- $\varepsilon(\mathbf{s}, t) \sim N(0, \sigma_\varepsilon^2)$  is the measurement error, time invariant (constant across time)
- **Example:** we have only one covariate altitude which is time invariant and  $y(\mathbf{s}, t)$  is the pollutant concentration. The global effect of the altitude is placed in  $\mathbf{x}_\beta(\mathbf{s}, t)$  which is constant across space and time but we have a **temporal correction** given by  $\mathbf{x}_z(\mathbf{s})' \mathbf{z}(t)$ . This because the effect of the altitude is different in winter and in summer. So if we put this covariate in both  $\mathbf{x}$ 's the  $\mathbf{x}_z(\mathbf{s})' \mathbf{z}(t)$  is an **adjustment** in respect of the global term  $\mathbf{x}_\beta(\mathbf{s}, t)$ . We can choose in our model where to place this covariate maybe just in  $\mathbf{x}_z(\mathbf{s})' \mathbf{z}(t)$  anyway we can do every combination we want. Best way to choose where to place the covariates is by → **cross-validation** and pick up the best model.

Insights:

- $y(\mathbf{s}, t)$ , where  $\mathbf{s}$  is a vector that indicates the spatial location such as (*latitude, longitude*) and  $t$  is the temporal dimension that will be discrete
- $\mathbf{z}(t)$  latent variable which is only temporal that means that is constant across space
- $\mathbf{x}_z(\mathbf{s})$  is a covariate vector, if we want that  $\mathbf{z}(t)$  interacts with some covariates, these covariates must be spatial (means that those don't change across time such as elevation while wind is spatial and temporal covariate). If the covariates are spatial and temporal we lose **model identifiability** or in other word we are **overfitting**.
- $\alpha w(\mathbf{s}, t) \sim GP$  and we have a  $GP$  for each time and they are independent → **assumptions/forcing**. Also we have a different  $\theta$  for each which can be  $\theta(t)$ , but since this parameter is not very well identifiable we can fix this parameter
- **Markovian dynamics** refers to the temporal evolution of the latent variable  $\mathbf{z}(t)$ . Specifically, it means that the state of the system at time  $t$  (represented by  $\mathbf{z}(t)$ ) depends only on the state at the previous time step,  $t-1$ , and not on any earlier time steps. So the current state retains information from the immediate past but is not influenced by events or states further back in time.
- $\mathbf{G}$  is a stable matrix which means that  $\mathbf{z}(t)$  doesn't diverge when time increases
- $\varepsilon(\mathbf{s}, t)$  scalar because  $y(\mathbf{s}, t)$  is scalar

### 4.2 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
- So using only the covariates can't describe this dynamics therefore we model the **persistency** using  $\mathbf{z}(t)$ . Even if we use one spatial model everyday we can't capture this effect.

### 4.3 Data matrix

- $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_T)$  is the data matrix
- Each  $\mathbf{y}_t$  is the vector of spatial observation at time  $t = 1, \dots, T$
- In each  $\mathbf{y}_t$ , missing data are possible
- Covariates are in a data array  $\mathbf{X} = \{\mathbf{X}_1, \dots, \mathbf{X}_T\}$ , each  $\mathbf{X}_t$  is a  $n \times b$  matrix where  $b$  is the number of covariates.  $\mathbf{X}$  cannot have missing data. One way to avoid missing data is through **interpolation/not kriging**
- So every  $t$ , based on the frequency of our dataset which is constant, we have a vector of observe quantity.
- Spatial locations  $\mathbf{s}$  can be fixed in our dataset but its common that maybe stations can be broken in some days and record from other stations. So generally in a more general scenario spatial locations  $\mathbf{s}$  can change. (e.g. smarthphone network)

### 4.4 Likelihood function

$$L(\Psi; \mathbf{Y}, \mathbf{W}, \mathbf{Z}, \mathbf{X}) = L(\Psi; \mathbf{Y} | \mathbf{W}, \mathbf{Z}, \mathbf{X}) \cdot L(\Psi; \mathbf{Z}) \cdot L(\Psi; \mathbf{W})$$

- The parameter vector is  $\Psi = \{\boldsymbol{\beta}, \alpha, \sigma_\varepsilon^2, \boldsymbol{\theta}, \mathbf{G}, \boldsymbol{\Sigma}_\eta\}$
- $\mathbf{W}$  and  $\mathbf{Z}$  have different dimensions.  $\mathbf{w}$  has the dimension of the network instead the dimension of  $\mathbf{z}$  is based on the number of covariante
- $\mathbf{W} = \{\mathbf{w}_1, \dots, \mathbf{w}_T\}$  for each time step we have a vector  $n \times 1 \Rightarrow$  matrix. In the model  $w$  is a scalar because we are explaining the output at spatial location  $\mathbf{s}$ , but we have  $n$  spatial locations ( $n \times 1$ )
- $\mathbf{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_T\}$  for each time step we have a vector  $q \times 1 \Rightarrow$  matrix
- $\mathbf{X}$  are covariates
- We can do the factorization of the Likelihood above because  $\mathbf{W} \perp \mathbf{Z} \rightarrow$  **model assumption**. And also in this factorization We are supposing to know  $\mathbf{W}$  and  $\mathbf{Z}$  but in practise we dont.

### 4.5 Log-likelihood function

$$\begin{aligned} -2\log(L_\Psi) &= T\log|\boldsymbol{\Sigma}_\varepsilon| + \sum_{t=1}^T \mathbf{e}_t' \boldsymbol{\Sigma}_\varepsilon^{-1} \mathbf{e}_t \\ &\quad + T\log|\boldsymbol{\Sigma}_\eta| + \sum_{t=1}^T (\mathbf{z}_t - \mathbf{G}\mathbf{z}_{t-1})' \boldsymbol{\Sigma}_\eta^{-1} (\mathbf{z}_t - \mathbf{G}\mathbf{z}_{t-1}) \\ &\quad + T\log|\boldsymbol{\Sigma}_w| + \sum_{t=1}^T \mathbf{w}_t' \boldsymbol{\Sigma}_w^{-1} \mathbf{w}_t \end{aligned}$$

*Innovation:*

$$\begin{aligned} \mathbf{z}(t) &= \mathbf{G}\mathbf{z}(t-1) + \boldsymbol{\eta}(t) \\ \boldsymbol{\eta}(t) &= \mathbf{z}(t) - \mathbf{G}\mathbf{z}(t-1) \end{aligned}$$

Where:

- $\mathbf{e}_t = \mathbf{y}_t - \mathbf{X}_{\beta,t}\boldsymbol{\beta} - \mathbf{X}_{z,t}\mathbf{z}_t - \alpha\mathbf{w}_t$
- $\boldsymbol{\Sigma}_\varepsilon = \sigma_\varepsilon^2 \mathbf{I}_n$ , with  $\mathbf{I}_n$  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)
- Maximizing the log likelihood means find the  $\mathbf{G}$  and the  $\boldsymbol{\Sigma}_\eta$  which explains the innovation that we are observing  $\Rightarrow$  but we are not observing  $\mathbf{z}$  such as  $\mathbf{w}$

## 4.6 Model estimation – EM algorithm

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

## 4.7 DCM alternative formulas

$$y(\mathbf{s}, t) = \mathbf{x}_\beta(\mathbf{s}, t)' \boldsymbol{\beta} + \mathbf{x}_z(\mathbf{s})' \mathbf{z}(t) + w(\mathbf{s}, t) + \varepsilon(\mathbf{s}, t)$$

$$\mathbf{z}(t) = \mathbf{G}\mathbf{z}(t-1) + \boldsymbol{\eta}(t)$$

$$w(\mathbf{s}, t) = \sum_{j=1}^c a_j x_j(\mathbf{s}, t) w_j(\mathbf{s}, t)$$

$$y(\mathbf{s}, t) = \mu(\mathbf{s}, t) + w(\mathbf{s}, t) + \varepsilon(\mathbf{s}, t)$$

$$\mu(\mathbf{s}, t) = x_\beta(\mathbf{s}, t) \boldsymbol{\beta} \quad \text{fixed effect model}$$

$$w(\mathbf{s}, t) = \sum_{j=1}^c a_j x_j(\mathbf{s}, t) w_j(\mathbf{s}, t) + x_z(\mathbf{s}, t) \mathbf{z}(t) \quad \text{random effect model}$$

$$\mathbf{z}(t) = \mathbf{G}\mathbf{z}(t-1) + \boldsymbol{\eta}(t) \quad \text{Markovian dynamics}$$

## 5 Lesson 5

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

$$y(\mathbf{s}, t) = \mathbf{x}_\beta(\mathbf{s}, t)' \boldsymbol{\beta} + \alpha z(\mathbf{s}, t) + \varepsilon(\mathbf{s}, t)$$

$$z(\mathbf{s}, t) = g z(\mathbf{s}, t - 1) + \eta(\mathbf{s}, t)$$

Remarks:

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

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

$$\mathbf{y}(\mathbf{s}, t) = \mathbf{X}_\beta(\mathbf{s}, t)' \boldsymbol{\beta} + \mathbf{z}(\mathbf{s}, t) + \boldsymbol{\varepsilon}(\mathbf{s}, t)$$

$$\mathbf{z}(\mathbf{s}, t) = \mathbf{G} \mathbf{z}(\mathbf{s}, t - 1) + \boldsymbol{\eta}(\mathbf{s}, t)$$

Remarks:

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

Insights and differences in respect to DCM:

- we have only 1 latent variable which is  $z(\mathbf{s}, t)$  and this a map.
- With the DCM you tend to overfit for each time step u are able to fit something which is spatial and which is independent from the previous time step and u have no constraints. U can change  $w(\mathbf{s}, t)$  at the previous time step.
- DCM has high  $R^2$  but it doesn't perform really well in cross validation because u are overfitting

### 5.3 Model estimation

- The HDGM is estimated similarly to the DCM
- But we only have the  $z(\mathbf{s}, t)$  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).

## 6 Lesson 6

### 6.1 Towards spatio-temporal functional models

Remarks:

- 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. In 3D space, one dimension may be sampled at higher resolution than the others

Disadvantages:

- Usually the original data set is very large and so the computational burden
- 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)

### 6.2 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

### 6.3 Basis Functions

We are going to choose B-spline bases as it is the most flexible basis to describe complex functions shape. We could also use other kind of basis, but one thing to keep in mind is that the number of basis function must be fixed before model estimation. High  $p$  (number of basis) usually means better  $R^2$ , but overfitting may be an issue. Usually when using an FDA approach it is recommended to use a high  $p$ , but in the case of f-HDGM this is not viable since the covariance matrices involved have dimension  $(n * p) \times (n * p)$ . Since  $n$  is usually large, a high  $p$  would make the estimation unfeasible.

Remarks:

- 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

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

- We now want to model the generic observation  $y(\mathbf{s}, t, h)$  where:
  - $\mathbf{s}$  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  $s$  describes the generic location across the globe

## 6.5 The functional HDG model in D-STEM

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

Remarks:

- $\eta(\mathbf{s}, t) \sim GP(0, \rho(||\mathbf{s} - \mathbf{s}'||; \boldsymbol{\theta}))$  is correlated over space but IID over time.
- $z(\mathbf{s}, t)$  is not scalar anymore and has Markovian dynamic
- $G$  is a diagonal transition matrix with diagonal elements in the  $p \times 1$  vector  $g$
- $\varepsilon(\mathbf{s}, t) \sim N(0, \sigma_\varepsilon^2(h))$  is a Gaussian measurement error independent in space and time. This means that the variance changes for each point of  $h$ , but it isn't influenced by position or time.
- $\boldsymbol{\phi}(h)$  is a  $p \times 1$  vector of basis functions evaluated at  $h$ , while  $\mathbf{c}_\varepsilon$  is a vector of coefficients to be estimated.

$$\log(\sigma_\varepsilon^2(h)) = \boldsymbol{\phi}(h)' \mathbf{c}_\varepsilon$$

- The model parameter set is  $\Psi = (\mathbf{c}_\varepsilon', \mathbf{c}_\beta', \mathbf{g}', \mathbf{v}', \boldsymbol{\theta}')'$

Insights:

- Since  $\varepsilon$  depends on parameter  $h$  this means that we are assuming  $\varepsilon$  to be heteroskedastic across the domain of the  $h$  variable. This means that for each value of  $h$  we have a different value of  $\varepsilon$
- The overall idea now is to have basis function and find the optimal coefficients to

## 7 Approfondimenti

$$\hat{w}(\mathbf{s}) \in \mathbb{R}, \text{ latent variable}$$

$$w(\mathbf{s}) \sim GP(0, \rho(||\mathbf{s} - \mathbf{s}'||; \boldsymbol{\theta}))$$

### 7.1 Kriging

**kriging:** is a spatial interpolation and prediction technique used in geostatistics to estimate values at unobserved locations based on the values observed at nearby locations.

- kriging also known as Gaussian process regression, is a method of interpolation based on Gaussian process governed by prior covariances
- Under suitable assumptions of the prior, kriging gives the best linear unbiased prediction (BLUP) at unsampled locations
- it provides not only point predictions but also estimates of the prediction uncertainty.
- kriging can be used to predict the values of the latent variable,  $w(\mathbf{s})$ , at unsampled spatial locations. This is particularly useful when you want to estimate pollutant concentrations at locations not covered by monitoring stations.
- Kriging predicts the value of a function at a given point by computing a weighted average of the known values of the function in the neighborhood of the point. The method is closely related to regression analysis. Both theories derive a best linear unbiased estimator based on assumptions on covariances

Steps:

- The first step in spatiotemporal kriging is to estimate the spatiotemporal variogram. This involves calculating the variance of the differences in values at different spatial locations and time points, considering both spatial and temporal distances. The variogram is used to model the spatiotemporal correlation structure of the data.
- The variogram is related to the covariance function, which represents the strength of the correlation between observations at different locations and times. The covariance function is used to construct a covariance matrix that accounts for both spatial and temporal dependencies.
- A kriging system is set up based on the covariance matrix, and the system includes observations at both observed and unobserved locations and times. The kriging system is then solved to obtain kriging weights.
- The kriging weights obtained from the system are used to predict values at unobserved spatial locations and time points. The kriging prediction is a linear combination of observed values with weights determined by the kriging system.
- Similar to traditional kriging, spatiotemporal kriging provides estimates of the prediction uncertainty. The kriging variance at each unobserved location and time point reflects the uncertainty associated with the predictions.
- The performance of the spatiotemporal kriging model can be assessed through cross-validation or other validation techniques to ensure its reliability and generalizability.

Key points related to kriging:

- **Variogram:** Kriging relies on the variogram, which quantifies the spatial correlation or covariance between data points at different spatial lags. The variogram model is used to describe how the spatial correlation decreases as the distance between locations increases.
- **Spatial Weights:** Kriging assigns spatial weights to observed data points, and these weights are used to estimate values at unobserved locations. The weights depend on the spatial correlation structure defined by the variogram.

## 7.2 Gaussian process

- a Gaussian process is a **stochastic process** (a collection of random variables indexed by time or space), such that every finite collection of those random variables has a multivariate normal distribution, i.e. every finite linear combination of them is normally distributed.
- The distribution of a Gaussian process is the joint distribution of all those (infinitely many) random variables, and as such, it is a distribution over functions with a continuous domain (time or space)
- Gaussian processes can be seen as an infinite-dimensional generalization of multivariate normal distributions. Gaussian  $\sim$  Normal
- A Gaussian stochastic process is strongly stationary or strict-sense stationary if and only if it is weak-sense stationarity, wide-sense stationarity (WSS).
- Basic aspects that can be defined through the covariance function are the process' stationarity, isotropy, smoothness and periodicity
  - If the process is stationary, the covariance function depends only on  $\mathbf{x} - \mathbf{x}'$
  - If the process depends only on  $|\mathbf{x} - \mathbf{x}'|$  (E. distance) then the process is considered isotropic.
  - A process that is concurrently stationary and isotropic is considered to be homogeneous
- Gaussian processes translate as taking priors on functions and the smoothness of these priors can be induced by the covariance function. If we expect that for "near-by" input points  $\mathbf{x}$  and  $\mathbf{x}'$  their corresponding output points  $\mathbf{y}$  and  $\mathbf{y}'$  then the assumption of continuity is present. If we wish to allow for significant displacement then we might choose a rougher covariance function.

### Conclusions:

- $w(\mathbf{s}_i)$  is a random variable and given a specific  $\mathbf{s}_i$ , the value of  $w(\mathbf{s}_i) \in \mathbb{R}$  comes from a distribution
- 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
- We impose that  $w(\mathbf{s}) \sim GP$ , and then im going to take a realization of a gaussian process and use it as my latent variable per each point in my dataset. I don't have any troubles because the  $GP$  is countinous and is definded everywhere
- When we talk about  $w(\mathbf{s})$  we are referring to the whole realization so the value of each  $w(\mathbf{s})_i$

### Insights:

- a stationary process is a stochastic process whose unconditional joint probability distribution does not change when shifted in time. Consequently, parameters such as mean and variance also do not change over time.
- Since stationarity is an assumption underlying many statistical procedures used in time series analysis, non-stationary data are often transformed to become stationary. The most common cause of violation of stationarity is a trend in the mean, which can be due either to the presence of a unit root or of a deterministic trend.
- A trend stationary process is not strictly stationary, but can easily be transformed into a stationary process by removing the underlying trend, which is solely a function of time. Similarly, processes with one or more unit roots can be made stationary through differencing. An important type of non-stationary process that does not include a trend-like behavior is a cyclostationary process, which is a stochastic process that varies cyclically with time.

### 7.3 Apulia paper

- We applied the HDGM to model two pollutants + include in the same model some meteorological variables **observed in another network** also as output of the model
- So we have a model with 8 variables where 2 are pollutants and 6 meteorological variables
- Why didn't we use the meteorological variables as covariates? because the grids are different, so one station might measure the PM10 but not all the meteorological variables. If I would like to use the meteorological variables as covariates I have to interpolate the covariates using any algorithm for spatial interpolation, but the problem is everytime we interpolate we have the number of the interpolation but we don't have the uncertainty. Another thing is that we could have missing values measuring the meteorological variables → and they are not allowed.
- For us the covariates have 0 uncertainty because they are just numbers
- the  $\mathbf{V}$  correlation matrix, considers the cross-correlation between all 8 variables
- The estimated model once we ran the EM was to use it in order to do spatial interpolation
- We basically used the information from the meteorological station to predict the pollutants concentration in the area where the station is located
- **The problem** is that when we consider in the model pollutants + meteorological variables we **lose the causality** of the effect that the weather has on the pollutants concentration, usually it's the weather that affects the pollutants concentration and not the opposite but our model doesn't know this. So this is the limit of this model, some outputs may not make sense.

### 7.4 Project clarifications

- **Standardizzare** significa trasformare una distribuzione di dati in modo che abbia una media di zero e una deviazione standard di uno. Questo processo è comunemente utilizzato nelle analisi statistiche e nei modelli di machine learning per rendere più facile la comparazione tra variabili che potrebbero avere scale di misura diverse o per facilitare l'interpretazione dei risultati.

La standardizzazione di una variabile  $X$  può essere espressa matematicamente attraverso la seguente formula:

$$X \sim d(\mu, \sigma^2) \rightarrow Z \sim d(0, 1)$$

$$Z = \frac{X - \mu}{\sigma}, \quad -1 \leq Z \leq 1$$

$X$	<i>è il valore originale della variabile</i>
$Z$	<i>è il valore standardizzato della variabile <math>X</math></i>
$\mu$	<i>è la media della distribuzione dei dati</i>
$\sigma$	<i>è la deviazione standard della distribuzione dei dati</i>

- La radice dell'errore quadratico medio (**RMSE**) è una misura della discrepanza tra i valori osservati e i valori predetti da un modello statistico o di previsione. È comunemente utilizzato per valutare la bontà di adattamento di un modello

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

Standardizzando, la mia response variable è compresa tra -1 e 1, dunque gli errori che commetto sono in questa scala. I residui sono poi elevati al quadrato e messi sotto radice, posso dunque vedere questo come un errore in percentuale che commetto.

- La **normalizzazione** è che è un tipo di standardizzazione in cui i dati vengono trasformati in modo che cadano entro un intervallo specifico, spesso tra -1 e 1 o tra 0 e 1
- Il **coefficiente di determinazione**, comunemente indicato come  $R^2$  (R al quadrato), è una misura statistica che indica quanto bene un modello si adatta ai dati osservati. In altre parole,  $R^2$  misura la proporzione di varianza della variabile dipendente che viene spiegata dalle variabili indipendenti nel modello.

**Sum of squares total**

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2 = SSR + SSE$$

**Sum of squares Error**

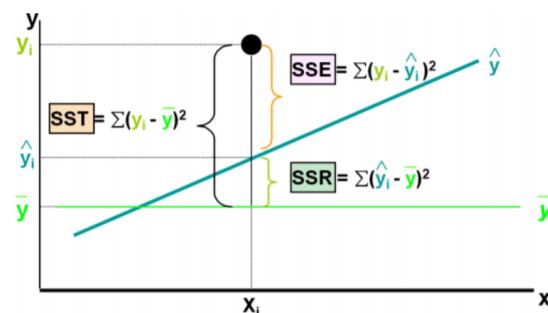
$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n \varepsilon_i^2$$

**Sum of squares due to regression (SSR)**

$$SSE = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

**R squared**

$$R^2 = 1 - \frac{SSE}{SST}$$



Se ottieni un coefficiente di determinazione  $R^2$  **negativo**, ciò indica che il tuo modello di regressione è peggio del semplice modello medio. In altre parole, il modello non spiega alcuna variazione o addirittura peggiora la capacità di predizione rispetto a una previsione basata semplicemente sulla media della variabile dipendente.

- **LBQ Test (Ljung-Box Q-test):** Il test LBQ viene utilizzato per verificare l'autocorrelazione residua nei residui di un modello di serie storica. L'**ipotesi nulla** del test LBQ è che non ci sia autocorrelazione nei residui. Se il risultato del test è 0, ciò indica che non vi è autocorrelazione residua nei residui del modello. Se il risultato è diverso da 0, indica che vi è autocorrelazione residua nei residui del modello.

**H0:** no autocorrelazione tra residui

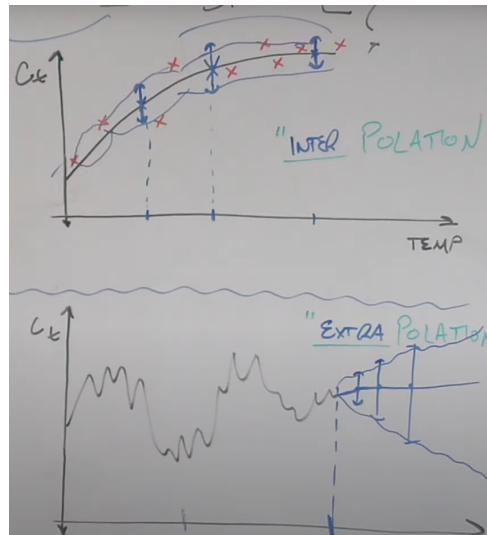
- **ARCH Test (Autoregressive Conditional Heteroskedasticity):** Il test ARCH viene utilizzato per verificare la presenza di eteroschedasticità condizionale nei residui di un modello di serie storica. L'**ipotesi nulla** del test ARCH è che non ci sia eteroschedasticità condizionale nei residui. Se il risultato del test è 0, ciò indica che non vi è eteroschedasticità condizionale nei residui del modello. Se il risultato è diverso da 0, indica che vi è eteroschedasticità condizionale nei residui del modello.

**H0:** residui omoschedastici

## 8 Time series analysis

### 8.1 Introduction:

- **Interpolation:** predictions from within a range of observable data. Based on the number of observable data  $u$  have we are going to build a more accurate model and residuals that are retrieved from the model are pretty much the same in terms of value for each data point. This happens because the uncertainty is the same within the data and becomes higher when we are going outside from the observable data → non time series problems are related to interpolation problems.
- **Extrapolation:** the further we go from our data points the higher the uncertainty is on our prediction. If I'm making a two days prediction in the future, the uncertainty is going to propagate (summable), so the uncertainty of the prediction is going to be added to the uncertainty of yesterday → Time series problems are related to extrapolation problems.



### 8.2 Autocorrelation (ACF) and partialcorrelation (PACF) function:

- **ACF:** tell us the correlation the effect at ago lags and the actual one, taking into account the direct and indirect correlation effect
- **PACF:** we only care of the direct effect, we don't care of the effect of indirect correlations

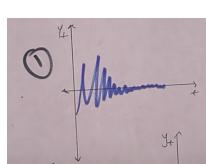
$$J_{s(t-2)} \rightarrow J_{s(t-1)} \rightarrow J_{s(t)}$$

$\text{corr}(J, M) = \text{direct correlation effect } J \rightarrow M$

$\text{corr}(J, M) = \text{indirect correlation effect } J \rightarrow F \rightarrow M$

### 8.3 Stationarity

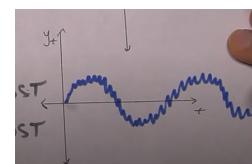
- Most of the models assume that the time series that we are trying to use them on is stationary
- **Def:** a time series is stationary when:  $\mu$  is constant over time,  $\sigma$  is constant over time and there is no seasonality
- Seasonality periodic behaviour over time that is predictable



$\mu$  is constant  
no seasonality  
 $\sigma$  is not constant



$\mu$  is not constant (shifting over time)  
no seasonality  
 $\sigma$  is constant (no fluctuations)



$\mu$  is constant  
seasonality (periodic component)  
 $\sigma$  is constant

- white noise signal is stationary by definition
- how to check if the process is stationary: visually, global versus local tests (divide the process into chunks), statistical tests such as "Augmented dickey-fuller (ADF test)"
- How to make a time series stationary:

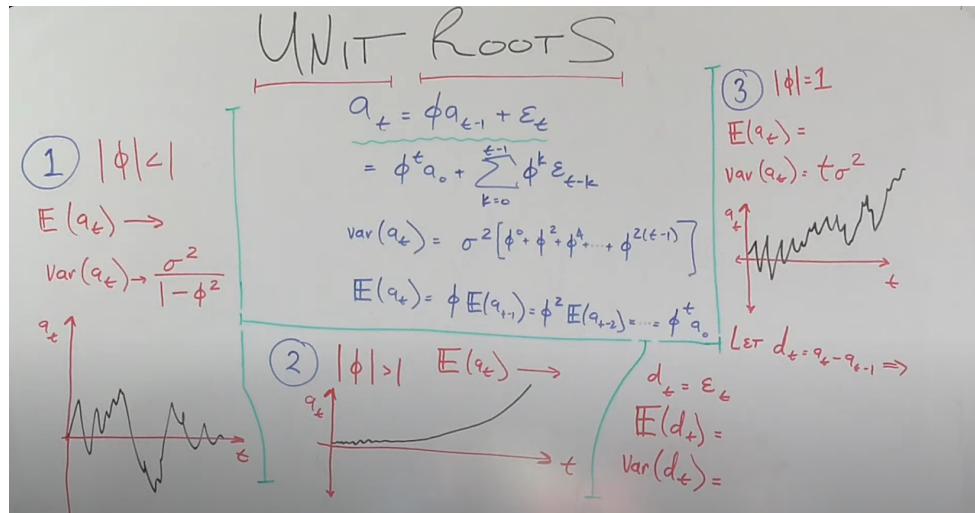
*y(t) is not stationary*

$$y_t = \beta_0 + \beta_1 t + \varepsilon_t$$

$$\begin{aligned} z_t &= y_t - y_{t-1} \\ &= (\beta_0 + \beta_1 t + \varepsilon_t) - (\beta_0 + \beta_1(t-1) + \varepsilon_{t-1}) \\ &= \beta_1 + \varepsilon_t - \varepsilon_{t-1} \\ \mathbb{E}[z_t] &= \beta_1 (\text{constant}) \\ \text{var}[z_t] &= 0 + k^2 + k^2 = 2k^2 (\text{constant}) \end{aligned}$$

## 8.4 Unit Roots

- If we have a time series with a unit root then the series is not stationary, so we have to apply some transformation in order to get it stationary
- We have unit root in example 3:  $|\phi| = 1$
- AR(1) it can be written as a  $MA(\infty)$ , try to write recursively  $AR(1)$  starting from  $a_t$
- the first value of the time series:  $a_0$



*Example 1:  $|\phi| < 1$*

$$\phi = 0.5$$

$$\mathbb{E}[a_t] = 0 \rightarrow \phi^\infty = 0$$

$$\text{var}[a_t] = \frac{\sigma^2}{1 - \phi^2}$$

Stationary ts

*Example 2:  $|\phi| > 1$*

$$\phi = 2$$

$$\mathbb{E}[a_t] = 0 \rightarrow \phi^\infty = \pm\infty$$

Non stationary ts

*Example 3:  $|\phi| = 1$*

$$\mathbb{E}[a_t] = a_0$$

$$\text{var}[a_t] = t \cdot \sigma^2$$

Non stationary ts but  
first difference

*First difference*

*new series:  $d_t = a_t - a_{t-1} = \varepsilon_t$*

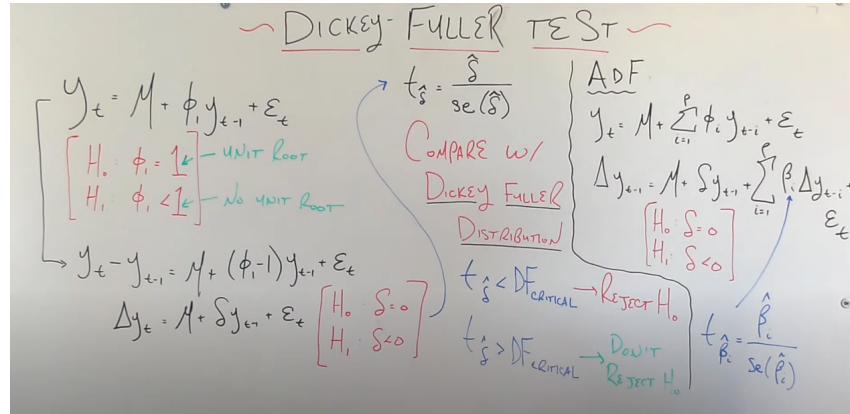
$$\mathbb{E}[d_t] = 0$$

$$\text{var}[d_t] = \sigma^2$$

*Now is stationary*

## 8.5 Dickey Fuller Test and Augmented Dickey Fuller Test

- Test in order to check if the time series is stationary (mean, variance are constant over time and no seasonality)
- Dickey Fuller Test assumes that our time series in question is a  $AR(1)$  process
- $se$ : standard error
- We compare the t-statistic against the dickey fuller distribution and not the normal distribution
- We extend this process to a more complex models than  $AR(1)$  and that's how **Augmented Dickey Fuller Test** comes into play. We also can calculate the t-distribution for each  $\beta_i$  and we can see if those  $\beta_i$  are significative or not



## 8.6 White noise

- white noise is a type of time series which has:

$$mean = 0$$

st. dev is constant with time

$$corr. between lags is 0$$

WN is not predictable

- Whenever we are doing a ts analysis on some  $y_t$  we always assume that  $y_t = signal + noise$ , where I can predict the *signal* and *noise*  $\sim WN$ , so unpredictable. If we are doing the things right we will end up with *WN* residuals and we can tell for sure our estimated model is really good.
- To test this we can: visual tests, check ACF and global vs local checks have to match (window that I shift every time and I calculate mean and variance and those have to be the same each time).
- The *ACF* checks for correlation between lags, if we see anything outside the bands we have a correaltion in that point

## 8.7 Backshift operator: lag operator

- Write in a more compact way the series
- we are going to trasform the first equation
- we bring all the *AR* terms on one side and the *MA* on the other one
- **lag operator:**  $L$  or  $B$

$$Ly_t = y_{t-1}, \quad L^2y_t = y_{t-2}$$

$$\begin{aligned} ARMA(3,3) : \\ Y_t &= \phi_1 Y_{t-1} + \phi_2 Y_{t-2} + \phi_3 Y_{t-3} + \\ &\quad \theta_1 \epsilon_{t-1} + \theta_2 \epsilon_{t-2} + \theta_3 \epsilon_{t-3} + \epsilon_t \\ \Rightarrow Y_t - \phi_1 Y_{t-1} - \phi_2 Y_{t-2} - \phi_3 Y_{t-3} &= \\ \epsilon_t + \theta_1 \epsilon_{t-1} + \theta_2 \epsilon_{t-2} + \theta_3 \epsilon_{t-3} &\Rightarrow \end{aligned}$$

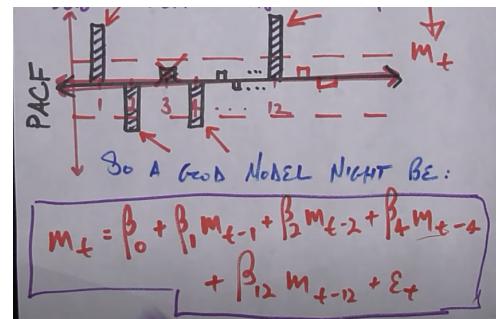
$$y_t - \phi_1 Ly_t - \phi_2 L^2 y_t - \phi_3 L^3 y_t = \epsilon_t - \theta_1 L \epsilon_t - \theta_2 L^2 \epsilon_t - \theta_3 L^3 \epsilon_t$$

$$(1 - \phi_1 L - \phi_2 L^2 - \phi_3 L^3)y_t = (1 - \theta_1 L - \theta_2 L^2 - \theta_3 L^3)\epsilon_t$$

$$\Phi(L)y_t = H(L)\epsilon_t$$

## 8.8 Autoregressive Model

- Autoregressive that means its a regression, so we are trying to predict something based on past values of the same thing
- In statistics if there are 2 models that have the same behavior we are going to prefer the simplest one (Occam's razor)
- notations for the exercise:  $m_t$  quantity of milk per month,  $m_{t-12}$  quantity of milk last year
- We want to use only the significative lags in order to predict the quantity of milk, if Im going to use all of them our model will be too tuned and we will get overfitting
- **PACF measure the direct correlation** and not taking into the account the indirect ones ( $m_{t-3} \rightarrow m_t$  removing  $m_{t-2}$  and  $m_{t-1}$ ). So we are going to use into the model the band that are statistically different than 0



## 8.9 Evaluate time series model

- Run machine learning: **train test split**, we train our model on 70% of the data and test on 30%
- Whenever we have a pattern in the residuals there is some dynamics that we didnt caputure with our model and so there is variability left inside the residuals → not great sign
- We are less and less certain on the predictions we are making moving further into the future and we arrive to predict essentially the mean of the time series. So splitting the data into training and testing we will see that the predictions will be okay initially and drop to mean widening the window of the prediction
- we can calculate the **mean absolute percent error**
- **Rolling forecast origin**: predict one month in advance each time and then average all the predictions. We need to fit the model every single time

**Train on months:**  $1, 2, \dots, k-3$

**Predict month**  $k-2$

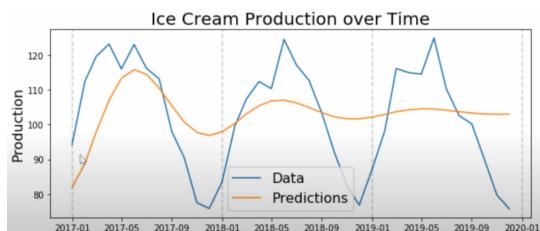
**Train on months:**  $1, 2, \dots, k-3, k-2$

**Predict month**  $k-1$

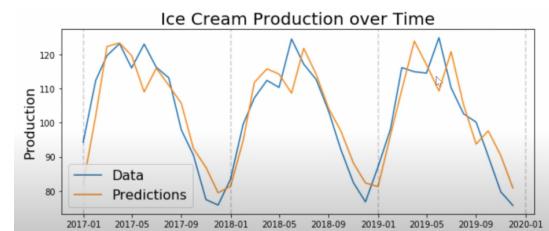
**Train on months:**  $1, 2, \dots, k-3, k-2, k-1$

**Predict month**  $k$

*average all the predictions*



splitting dataset into train and test



Rolling forecast origin

## 9 Moving average model

- $\hat{f}_t$  predicted number of cupcakes
- $\varepsilon_t$  error the professor tells u
- $f_t$  how many cupcakes u should have bought that month
- this trend seems to be centered in 10. This is why this model is called moving average, because we have our average of 10 but that average its moving all about that average but still staying centered sort of over there → always center at 10
- The model we are seeing its a MA(1) so we are taking into account one error from the last month
- In order to estimate the model from the data check ACF function and see the correlated lags
- $\mu$  is a additive constant, in our IMAD models is 0

$$MA(1) : \hat{f}_t = \mu + \phi_1 \varepsilon_{t-1} \quad \varepsilon_t \sim N(\mu_\varepsilon, \sigma_\varepsilon^2) = (0, 1)$$

$$\mathbb{E}[\hat{f}_t] = \mathbb{E}[\mu + \phi_1 \varepsilon_{t-1}] = 10 + 0$$

$$MA(2) : f_t = \mu + \phi_1 \varepsilon_{t-1} + \phi_2 \varepsilon_{t-2} + \varepsilon_t$$

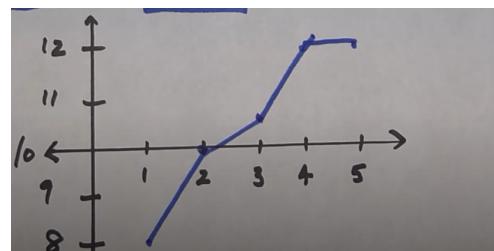
$$\hat{f}_t = \mu + \phi_1 \varepsilon_{t-1} + \phi_2 \varepsilon_{t-2}$$

Moving Average Model

$\hat{f}_t = \mu + \phi_1 \varepsilon_{t-1}$

$\mu = 10, \phi_1 = .5$

t	$\hat{f}_t$	$\varepsilon_t$	$f_t$
1	10	-2	8
2	9	1	10
3	10.5	0	10.5
4	10	2	12
5	11	1	12



### 9.1 Moving Average and ACF(Auto correlation function)

- In the correlation fuction the only terms that will be not equal to 0, will be the variance  $\mathbb{E}[\varepsilon_t^2]$  this will happen if we have overlapping lags
- $\mathbb{E}[\varepsilon_t \cdot \varepsilon_{t-k}] = 0, k \neq 0$ , because the errors are independent

$$ACF(k) = \text{corr}(x_t, x_{t-k}) = \mathbb{E}(x_t x_{t-k}) - \mathbb{E}(x_t) \mathbb{E}(x_{t-k})$$

Moving Average & ACF

MA(2):  $x_t = \mu + \phi_1 \varepsilon_{t-1} + \phi_2 \varepsilon_{t-2} + \dots + \phi_2 \varepsilon_{t-2} + \varepsilon_t$

$\text{corr}(x_t, x_{t-k}) = \mathbb{E}(x_t x_{t-k}) - \mathbb{E}(x_t) \mathbb{E}(x_{t-k})$

$\Rightarrow x_t : [\varepsilon_{t-1}, \varepsilon_{t-2}, \dots, \varepsilon_{t-2}, \varepsilon_t]$

$x_{t-k} : [\varepsilon_{t-k-1}, \varepsilon_{t-k-2}, \dots, \varepsilon_{t-k-2}, \varepsilon_{t-k}]$

$\neq 0 \text{ IFF } t-2 \leq t-k \Rightarrow k \leq 2$

$\neq 0 \text{ IFF } t-2 \leq t-k \Rightarrow k \leq 2$

So ACF looks like ...

$\text{var}(\varepsilon_t) = \mathbb{E}(\varepsilon_t^2) - \mathbb{E}(\varepsilon_t)^2$

## 9.2 Invertibility of Time Series: $MA(1) \Leftrightarrow AR(\infty)$

- Connection between AR and MA models
- $\frac{1}{1-\phi L}$  is an infinite geometric sum, where the common ratios is  $\phi L$ , the restriction is that the common ratios absolute value must me less than 1 if not the geometric series doesnt converge  $\rightarrow |\phi| < 1$  in this way we can say that the time series is invertible

$$\begin{array}{l}
 C_t \leftarrow \varepsilon_t \\
 \uparrow \\
 \varepsilon_{t-1} \leftarrow C_{t-1} \\
 \uparrow \\
 \varepsilon_{t-2} \leftarrow C_{t-2} \\
 \uparrow \\
 \vdots \\
 \varepsilon_{+3} \leftarrow C_{t-3} \\
 \uparrow \\
 \dots
 \end{array}
 \quad
 \left\{
 \begin{array}{l}
 C_t = -\phi \varepsilon_{t-1} + \varepsilon_t \\
 C_t = (1-\phi L) \varepsilon_t \\
 \frac{C_t}{1-\phi L} = \varepsilon_t \Rightarrow \\
 (1+\phi L + \phi^2 L^2 + \dots) C_t = \varepsilon_t \Rightarrow \\
 C_t + \phi C_{t-1} + \phi^2 C_{t-2} + \dots = \varepsilon_t \Rightarrow \\
 C_t = -\phi C_{t-1} - \phi^2 C_{t-2} - \dots + \varepsilon_t
 \end{array}
 \right.$$

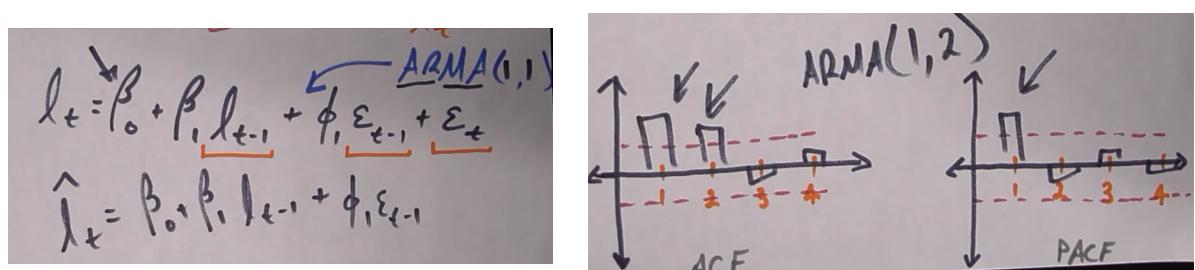
## 9.3 Invertibility of Time Series: $AR(1) \Leftrightarrow MA(\infty)$

- Same notes as before: important is that  $|\phi| < 1$
- so invertibility processes are not seprate entities but are the two sides of the same coin

$$\begin{array}{l}
 C_t \leftarrow \varepsilon_t \\
 \uparrow \\
 C_{t-1} \leftarrow \varepsilon_{t-1} \\
 \uparrow \\
 C_{t-2} \leftarrow \varepsilon_{t-2} \\
 \uparrow \\
 \vdots \\
 C_{t-3} \leftarrow \varepsilon_{t-3} \\
 \uparrow \\
 \dots
 \end{array}
 \quad
 \left\{
 \begin{array}{l}
 C_t = \phi C_{t-1} + \varepsilon_t \\
 (1-\phi L) C_t = \varepsilon_t \Rightarrow \\
 \frac{1}{1-\phi L} \varepsilon_t = C_t \\
 (1+\phi L + \phi^2 L^2 + \dots) \varepsilon_t = C_t \\
 \phi \varepsilon_{t-1} + \phi^2 \varepsilon_{t-2} + \phi^3 \varepsilon_{t-3} + \dots + \varepsilon_t = C_t
 \end{array}
 \right.$$

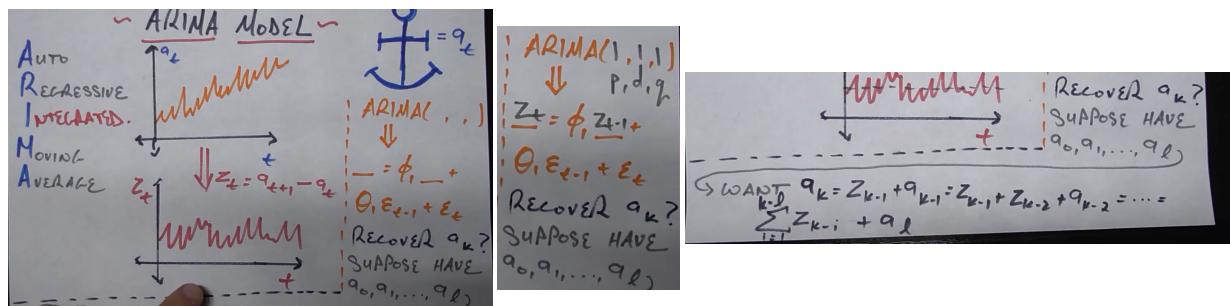
## 9.4 ARMA model

- $l_t$  : light bulbs demand at the current time, so how many should I build up
- $\beta_0$  : cosntant demand of light bulbs
- $\beta_1 l_{t-1}$  light bulbs needed to create last month
- $\phi_1 \varepsilon_{t-1}$  : error from the previous time period. Last month I made some prediction of how much light bulbs to create and my prediction was off something
- ACF helps to indentify the MA part and PACF the AR part of the ARMA



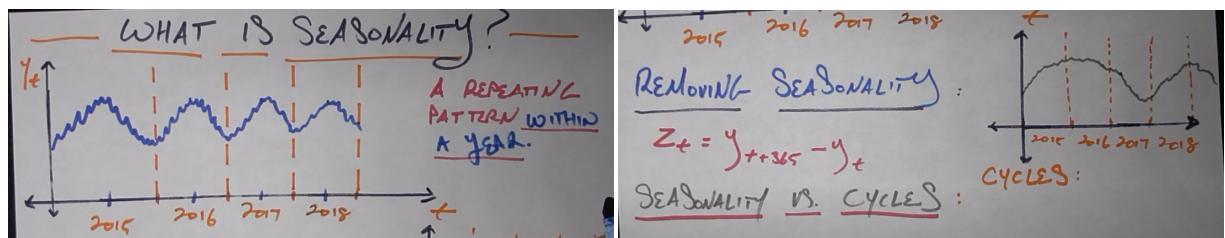
## 9.5 ARIMA model

- $q_t$  : number of anchor u sell every month
- we can't use ARMA model because the time series is not stationary because of a linear trend in the mean → in this case we can use ARIMA models
- **Integrated:** instead of predicting the time series itself we are going to predict the differences from one timestamp to the previous one. So we are basically creating a new time series which is  $z_t$
- Basically a function is going to a precise pattern and this pattern is constant, so if the function is going to increase +1 at  $t = 0$ , is going to do the same at  $t = 1$  and so one. The behavior of a function is constant over time as well as our time series, and that's why we have created the new time series  $z_t$
- Setting the parameter  $d = k$  we are going to differentiate it  $k$  times in a row getting a new time series. I'm going to stop the differentiate (and so the parameter  $d$ ) process once the time series is stationary. We want  $d$  to be as low as we can in order to build up a simple model
- What we want to predict is not the difference but the actual time series, so how we get back from  $z_t$  to  $q_t$ : ( $q_l$  last value we have data for)



## 9.6 What is Seasonality

- **Seasonality** is a repeating pattern within a year that is predictable, instead the **Cycles** is a pattern that shows up during the years and is an unpredictable component
- If a time series has a seasonal component it isn't stationary so we can't use any model we know so far.
- **Idea:** if I'm going to take in each period a value in the same position I'm expecting to get the same value (seasonality - repeating pattern). In order to remove it I'm going to delete the current value in the period with the previous one and get the innovation out of it



## 10 Roadmap time series

- Visualizing the time series
- distribution of the covariates
- Understand the imputation uncertainty
- response variable NH3 concentration
- use ACF and PACF on a subset pf the data (there is a function inside ACF that can skip the NA values)
- preliminary analisys on the missing data
- Cleaning dataset:
  - replace NAs with 0
  - replace missing values with average values
  - remove Nas values
  - missing robost method
- ARIMA doesnt have any problem with missing values, is a resistent estimation method on missing data without any assumption on them
- Model selection techniques: AIC, BIC, lasso, ridge
- use only the significative covariates and standardize it so will me in the same units
- Decompose the model into:

$$y(t) = \text{trend} + \text{seasonality} + \text{Cyclicity} + \text{innovation}$$

*seasonality: additive or multiplicative*

*detrending*

- Make the time series stationary and test the stationary
- Model the innovation using multiple models
- Validation of the dataset 70% and 30% and try other techniques, or block bootstrap (there is some link somewhere on the teams)
- Block bootstrap only if the data is stationary
- Residuals check  $e(t) \sim Wn$
- Compare models and choose the best one according to the Occam's razor
- try machine learning methods and deep learning ones and explain the theory behind ht method
  - Kaggle
  - LightGBM
  - Decision Trees
  - XGBoost
  - Long Short-Term Memory (LSTM)
  - Recurrent Neural Network (RNN)
  - Multi-Layer Perceptron (MLP)

*site – link*

- We can even delete the time component and deal with a regression model
- If we use auto arima we have to read the paper and see how it works
- Use RMSE in order to get the best model

## 11 Machine learning models

- **Ensemble learning:** combine multiple weak models/learners into one predictive model to reduce bias, variance and/or improve accuracy. These models are known as **weak learners**. The intuition is that when you combine several weak learners, they can become **strong learners**.
- **Gradient Boosting** is an ensemble learning algorithm used for both regression and classification tasks. Here's a brief description of how it works:
  - **Base Learners:** Gradient Boosting builds an ensemble of weak learners (usually decision trees) sequentially. Each weak learner tries to correct the errors made by the previous ones.
  - **Gradient Descent:** It optimizes a loss function by adding weak learners to the ensemble in a forward stage-wise manner. At each stage, the algorithm calculates the gradient of the loss function with respect to the ensemble's prediction, and fits a weak learner to the gradient.
  - **Boosting:** The subsequent weak learners focus on the mistakes (residuals) made by the previous ones. They are trained on the residuals of the previous learners rather than the original data.
  - **Weighted Voting:** The final prediction is made by a weighted sum of the predictions of all weak learners. The weights assigned to each learner depend on their individual performance in minimizing the loss function.
  - **Regularization:** Gradient Boosting employs regularization techniques to prevent overfitting, such as limiting the maximum depth of the trees, adding a shrinkage parameter to the predictions, and using subsampling of the data.
- **XGBoost:** An optimized distributed library for machine learning models in the gradient boosting framework, designed to be highly efficient, flexible, and portable. It features regularization parameters to penalize complex models, effective handling of sparse data for better performance, parallel computation, and more efficient memory usage

### 11.1 XGBoost

- XGBoost (eXtreme Gradient Boosting) is an implementation of the Gradient Boosting algorithm
- XGBoost is a decision tree-based ensemble algorithm that leverages the concept of gradient boosting to create a strong model
- **How it works:** works by creating a set of decision trees iteratively, each tree attempting to correct the mistakes of the previous tree. The algorithm employs a **gradient descent algorithm** to minimize a cost function, which is the sum of the errors of each tree in the ensemble.
- The final model is a weighted combination of all the decision trees, with each tree assigned a weight based on its contribution to the cost function.
- Steps:
  1. **Initialize the model:** the first step in XGBoost is to initialize the model by creating a base prediction. This prediction is usually the mean of the target variable for regression problems or the most common class for classification problems.
  2. **Fit the first tree:** to train the first tree of a model, features and residuals are used in a greedy manner where informative features are selected first. The residuals are the differences between the model's predictions and the actual values.
  3. **Compute the loss:** after training the first tree, the loss is computed by summing the errors made by all trees in the ensemble. The sample weights are then updated using the loss, with higher weights assigned to samples with higher errors.
  4. **Fit the next tree:** after computing the loss, the next step is to fit the next tree. This tree is trained on the updated weights and the residuals from the previous tree. The process is repeated until a predetermined number of trees is reached or the loss stops decreasing.
  5. **Make predictions:** The final step in XGBoost is to make predictions using the ensemble of trees. The prediction for each sample is the weighted sum of the predictions made by each tree in the ensemble.
- XGBoost utilizes several techniques to build its algorithm. These include:

- **Tree Building:** XGBoost uses a greedy algorithm to build decision trees. It starts with a single node and recursively splits the data into two child nodes based on the feature that maximizes the reduction in the loss function. The loss function is typically the negative log-likelihood for classification problems and the mean squared error for regression problems. The splitting process continues until a stopping criterion is met, such as a maximum depth or a minimum number of samples per leaf node. XGBoost also includes **tree pruning** to prevent overfitting, which removes nodes that do not contribute to the reduction in the loss function.
- **Gradient Descent:** gradient descent is an iterative method to find the minimum of a function by adjusting the parameters towards the negative gradient. In XGBoost, the loss function is the sum of the losses of each data point and a regularization term that penalizes large coefficients. The regularization term, either L1 or L2, makes the model simpler and prevents overfitting. interpretable.

## 11.2 Code:

In the context of XGBoost, the terms "Gain," "Cover," and "Frequency" represent different aspects of feature importance calculated by the algorithm:

**Gain:** This refers to the improvement in accuracy brought by a feature to the branches it appears in. Essentially, it measures the contribution of each feature to the model's performance. Higher gain values indicate that the feature is more important for making decisions in the model.

**Cover:** Cover stands for the relative quantity of observations concerned with a specific feature. It's the sum of second-order gradients for all instances in a tree that the feature is used to split. Higher cover values indicate that the feature is used to make more splits across the dataset.

**Frequency:** Frequency denotes the number of times a feature appears in all generated trees. It provides insights into how often a feature is used for splitting nodes across all the trees in the ensemble.

In summary, while Gain assesses the quality of a feature, Cover and Frequency offer information about the quantity and usage of the feature in the construction of the ensemble model. These metrics collectively help in understanding the importance and contribution of each feature towards the predictive performance of the model.

## 11.3 Understanding Facebook's Prophet

Prophet is a forecasting tool developed by Facebook's Core Data Science team. It's designed to make time series forecasting tasks more accessible to analysts and non-experts. Here's how it works:

1. **Automatic Forecasting:** Prophet is built to handle time series data with strong seasonal patterns and multiple seasonality. It automatically detects these patterns in the data and fits a model accordingly.
2. **Flexible Model:** Prophet uses a decomposable time series model with three main components: trend, seasonality, and holidays. The trend component models non-periodic changes over time, while the seasonality component captures periodic changes (daily, weekly, yearly). Holidays and special events can also be incorporated into the model.
3. **Bayesian Approach:** Prophet uses a Bayesian approach to model fitting, which allows for uncertainty estimation in the forecasts. This means that Prophet provides not only point forecasts but also uncertainty intervals around the predictions, which can be valuable for decision-making.
4. **Robustness to Missing Data and Outliers:** Prophet is designed to handle missing data and outliers in the time series gracefully. It employs a procedure called "imputation" to fill in missing values in the data, and it uses a robust method to identify outliers and mitigate their impact on the forecasts.
5. **User-Friendly Interface:** One of the key features of Prophet is its user-friendly interface. It provides a simple API that allows users to specify the time series data, as well as any additional information such as holidays or seasonality patterns. This makes it easy for analysts and data scientists to use Prophet even without extensive knowledge of time series forecasting techniques.
6. **Scalability:** Prophet is designed to be scalable and can handle large datasets efficiently. It's implemented in Python and is built on top of the Stan probabilistic programming language, which provides efficient algorithms for Bayesian inference.

<b>Max depth</b>	<b>min iterations</b>	<b>train RMSE</b>	<b>validation RMSE</b>
max_depth_1	998	6.088	7.325
max_depth_2	997	5.349	7.132
max_depth_3	745	4.552	7.097
max_depth_4	324	3.684	7.047
max_depth_5	336	2.758	6.943
max_depth_6	347	1.914	6.954

<b>Model</b>	<b>weekly</b>	<b>monthly</b>	<b>annual</b>
XGBoost	2.200e-16	2.200e-16	2.200e-16
Prophet	2.200e-16	2.200e-16	2.200e-16
LSTM	0.058	0.355	0.065

$$y(t) = g(t) + s(t) + h(t) + e(t)$$

<b>Model</b>	<b>train RMSE</b>	<b>validation RMSE</b>	<b>Shapiro test</b>	<b>Breusch-Pagan Test</b>
XGBoost	4.281	6.934	1.227e-08	0.186
Prophet	6.772	7.924	4.721e-09	0.240
LSTM	1.238	7.195	0.062	0.207e-03

<b>Epochs</b>	<b>train <i>RMSE</i></b>	<b>validation <i>RMSE</i></b>	<b>Shapiro test</b>	<b>Breusch-Pagan Test</b>
200	1.331	7.010	2.772e-05	0.253e-03
100	1.238	7.195	0.062	0.207e-03
50	2.223	7.173	4.409e-05	1.556e-05

<b>Model</b>	<b>weekly</b>	<b>monthly</b>	<b>annual</b>
200	0.025	0.543	0.299
100	0.058	0.355	0.065
50	2.116e-05	0.018	0.373