



Spatial Data Analysis for a Spatial World

Raquel Menezes

Centro de Matemática, Universidade do Minho
rmenezes@math.uminho.pt

Actividade INE - Meia Hora

June 25, 2025

Outline

① Introduction to Spatial Statistics

- ▷ Data Analysis and Visualization

② Spatial Data: **Areal, Point, Geostatistical**

- ▷ Understanding Spatial Correlation

③ Applications to Fisheries and Marine Sciences

- ▷ Try **Dynamic Maps**

The Beginning of Spatial Statistics

John Snow and the 1855 London Cholera Epidemic

Background

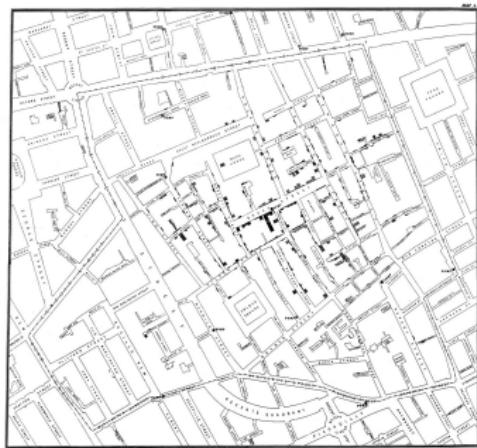
- **John Snow (1813-1858):** British physician well-known in epidemiology

The Cholera Epidemic

- **London 1854:** Severe cholera outbreak hits Soho, causing numerous deaths
- **Common Belief:** Cholera was believed to be spread by bad air

Snow's Insight

- **Hypothesis:** Snow theorized that **cholera was waterborne**, not airborne
- **Investigation:** He meticulously **mapped cholera cases in Soho**



The Breakthrough

- **Pump Analysis:** Snow identified a pattern centering around Broad Street water pump.
- **Action Taken:** Persuaded authorities to remove the pump handle, leading to a dramatic decline in cases.



Impact on Spatial Statistics

- **Mapping Cases:** Snow's use of maps to identify the outbreak's source is considered a foundational moment in spatial statistics.
- **Legacy:** Demonstrated the power of spatial data in understanding and solving public health issues.

Classic EDA and Its Main Objectives

We must consider:

- **Numerical summaries**, known as descriptive statistics (means, medians, quantiles, variance, ...)
- **Graphical summaries** related to preliminary analysis, through data visualization

Know your data!

- distributions (symmetric, normal, asymmetric)?
- data quality issues?
- outliers or extreme values?
- correlations and interrelations?
- subsets of interest?
- suggestions of functional relationships?

Sometimes, EDA (together with [data visualization](#)) can be the main objective!

Why Punctual Statistics Are Not Sufficient?

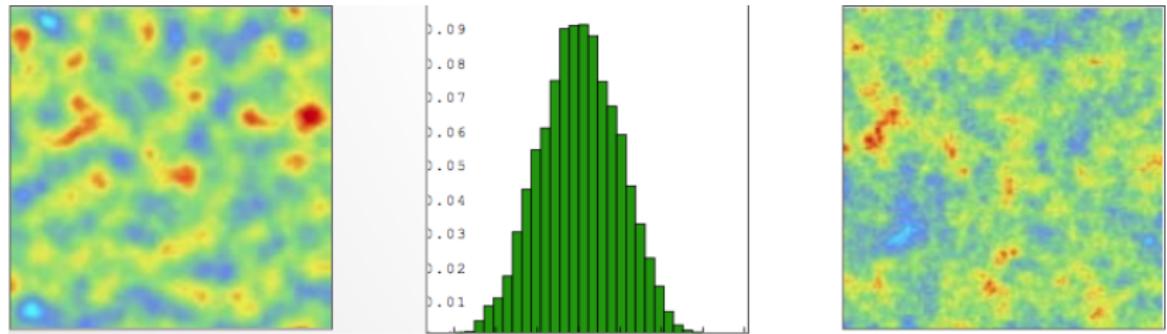


Figure: Two random fields (**highest values** and **lowest values**) with the same histogram.

An Example - Meuse River Data (*gstat* R package)

Top soil heavy metal concentrations (ppm), collected in a flood plain of river Meuse, near village Stein:

- cadmium (Ca)
- copper (Cu)
- lead (Pb)
- zinc (Zn)

| Data | Ca | Cu | Pb | Zn |
|-----------------|------|-------|--------|--------|
| Original | | | | |
| mean | 3.11 | 39.42 | 148.55 | 464.60 |
| median | 1.9 | 29.5 | 116.0 | 307.5 |
| st.dev | 3.49 | 23.40 | 110.18 | 376.57 |
| Logarit. | | | | |
| mean | 0.52 | 3.53 | 4.77 | 5.86 |
| median | 0.64 | 3.38 | 4.75 | 5.73 |
| st.dev | 1.21 | 0.51 | 0.67 | 0.73 |

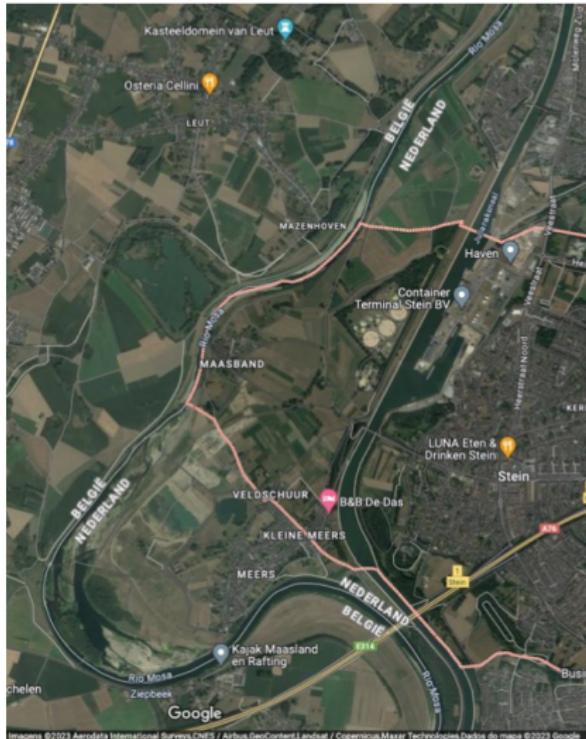
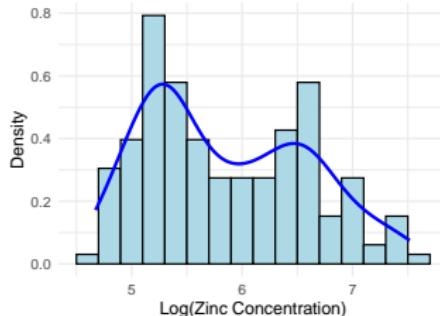


Table: Example of numerical summaries.

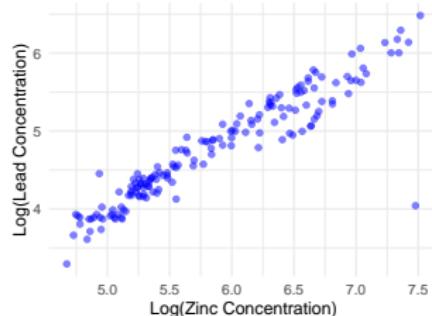
Examples of Graphical Summaries

Univariate, Bivariate and Multivariate

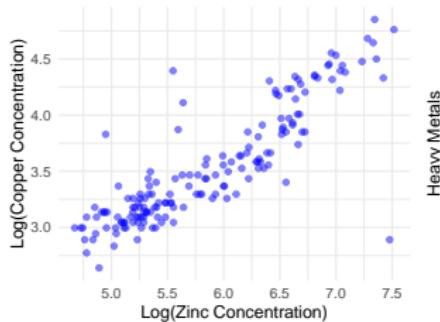
Histogram and Density Plot



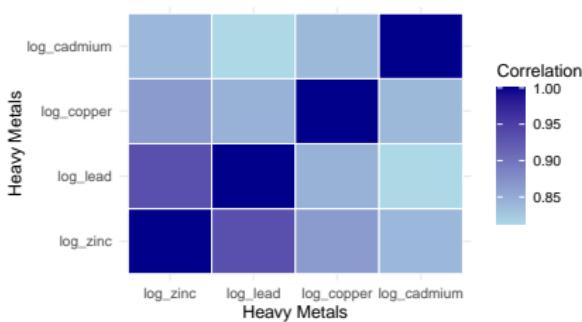
Bivariate Scatter Plot

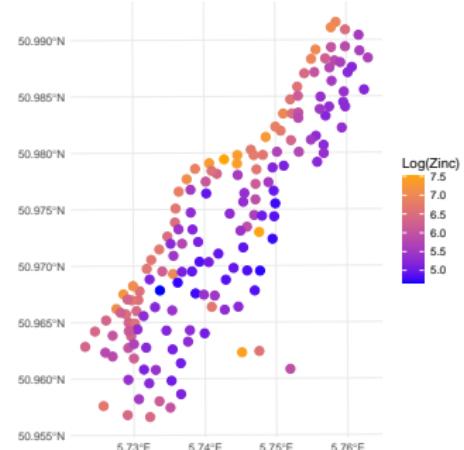
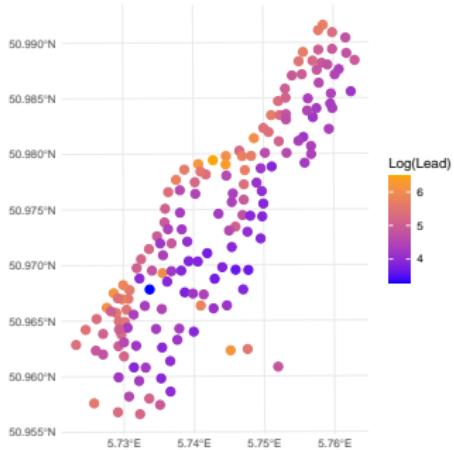
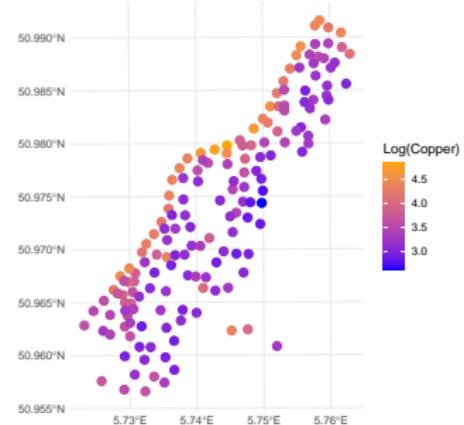
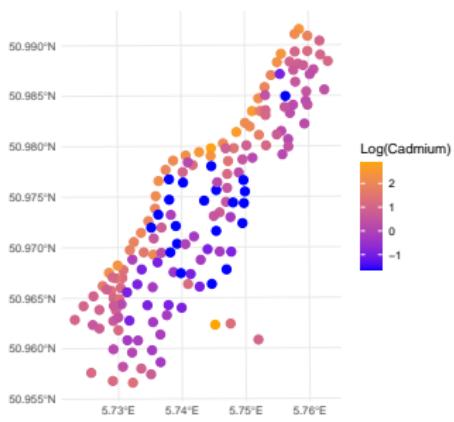


Bivariate Scatter Plot



Heatmap of Correlation Matrix





Summarizing remarks

Exploratory data analysis and data visualization:

- It helps to get a general sense of the data.
- You should always look at every variable - you will learn something!
- It's a **data-driven (model-free) approach**.

So it is **always well worth looking at your data, before moving on to any modelling approach!**

Outline

① Introduction to Spatial Statistics

- ▷ Data Analysis and Visualization

② Spatial Data: **Areal, Point, Geostatistical**

- ▷ Understanding Spatial Correlation

③ Applications to Fisheries and Marine Sciences

- ▷ Try DYNAMIC MAPS

Spatial Data

In spatial statistics, our process of interest $Y(x)$ is defined over a spatial region $x \in A$, and there are observations at specific locations x_1, \dots, x_n .

Depending on the **nature of the data** and the **spatial aggregation** we give them, we can differentiate three types of spatial data:

- ① Areal data, also known as **lattice data**
- ② Data referring to **point patterns**
- ③ Point-referenced data, known as **geostatistical data**

1. Lattice Data

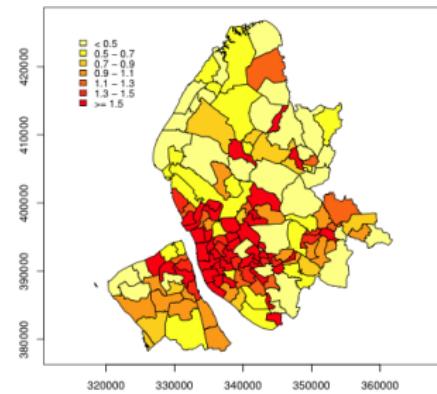
- Data referring to areas, or simply **lattice data**, represent an **aggregation of observations** $Y(x)$ over a predefined *unit of area*¹, such as
 - Number of cases of a specific disease in each area
 - Percentage of unemployment in each area
- The region of study is divided into a **finite collection of area units** with well-defined boundaries, such as
 - Autonomous communities of Spain
 - Districts or municipalities in Portugal
 - Countries in the European Union

In the modeling of lattice data, one must consider **whether adjacent regions have similarities**, in the sense that it is expected that nearby areas have more in common than distant areas.

¹Point x might represent the centroid of that unit area.

Example of Aggregated Data

Laryngeal Cancer, Northwest England (Bailey, 2008)



- Data on **876 cases of laryngeal cancer** diagnosed in 1982-1991, in **144 electoral wards** of Mersey and West Lancashire.
- Research allowed the definition of a **smoking prevalence indicator** in each district ('low', 'medium', or 'high').
- Available **annual air pollution measure** based on traffic flow.

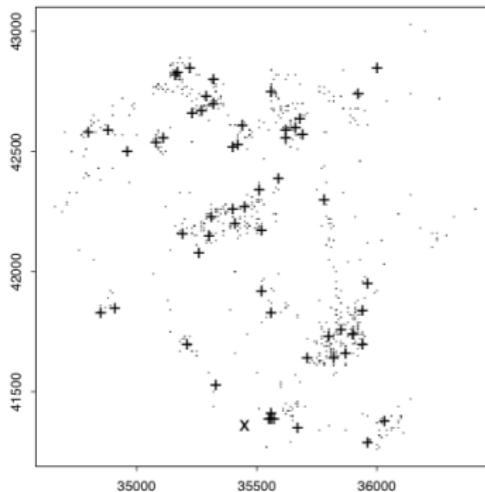
2. Data referring to Point Patterns

- The theory of point patterns emerged from the need to model the location of events of interest as random.
- The first applications were in the fields of ecology (e.g., locations where a rare bird is sighted) and forest sciences (e.g., location of a type of trees).
- However, their range of applications is very broad. We can consider the distribution of vessels at sea, or the distribution of known cases of a particular contagious disease.

In these various contexts, the goal is to study the spatial arrangement of the event of interest in space, to identify important areas (e.g., favorable areas for a given phenomenon).

Example of Point Patterns

Lung and Larynx Cancers (Diggle, Gatrell, and Lovett, 1990)



The map shows:

- cases of lung cancer (dots)
- cases of larynx cancer (crosses)
- a deactivated industrial incinerator (x symbol)

Important Questions:

- Do the cases show a *surprising* tendency to cluster?
- Does the disease risk vary spatially?
- Is the disease risk elevated near a specific location?

3. Geostatistical Data

- Geostatistical data (or point-referenced data) consist of a **stochastic process** $Y(x)$ associated with a **fixed set of locations x** over a **continuous spatial field $S(x)$** .
- The space is typically treated as **two-dimensional**, defined by its **longitude** and **latitude**, but it can also include altitude or depth to make it three-dimensional.
- Examples of geostatistical processes:
 - $S(x)$ represents the **pollution surface** over the city of Lisbon, and $Y(x)$ represents a pollution indicator measured at the monitoring station at location x .
 - In fishing, $S(x)$ may represent the **abundance surface** of a species and $Y(x)$ the catch of that fish at location x .

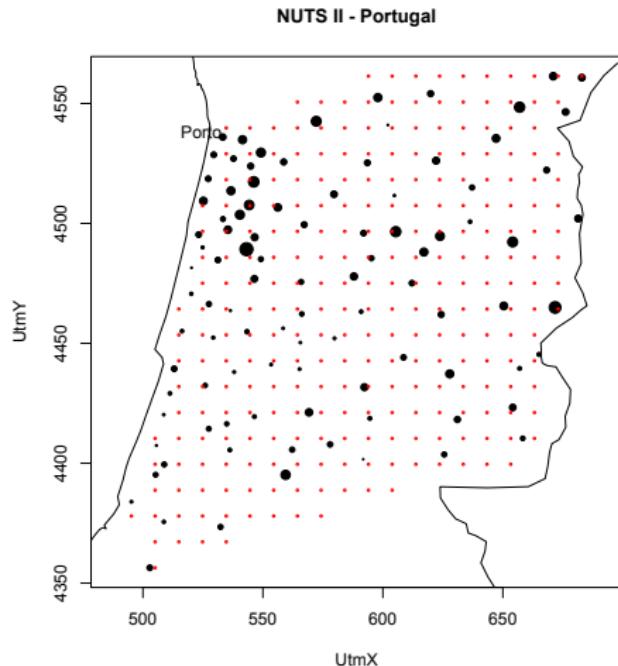
Modeling Geostatistical Data

- When modeling our **variable of interest**, conditioned on **predictors**, it is expected that the respective **residuals are spatially correlated**.
- Our main objective is to **infer the spatial correlation** underlying our data to improve prediction, using (for example) **kriging** techniques at unsampled locations.

Obviously, in all three sub-areas of Spatial Statistics, the **spatial dimension can be extended to the spatiotemporal domain** by adding the correlation of the variable of interest between temporal events (e.g., studying fish abundance every hour, day, etc).

Example of Geostatistical Data

Arsenic pollution data $Y(x)$ with $x \in A \subset I\!R^2$ (Garcia-Soidán and Menezes, 2017).



Note: Black circles identify the 98 measurements of As, with the circle diameter proportional to the observed value. Red points identify the grid points where As is to be estimated.

Fundamental Concept of “Spatial Dependence”

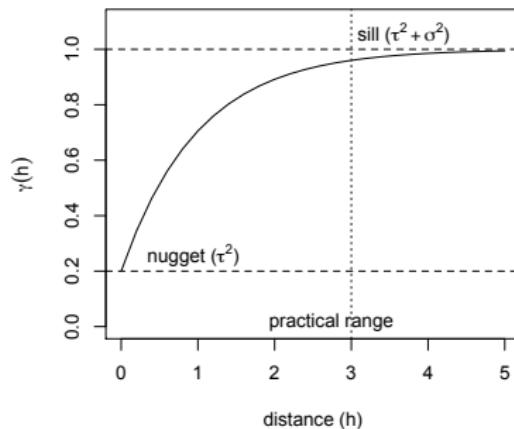
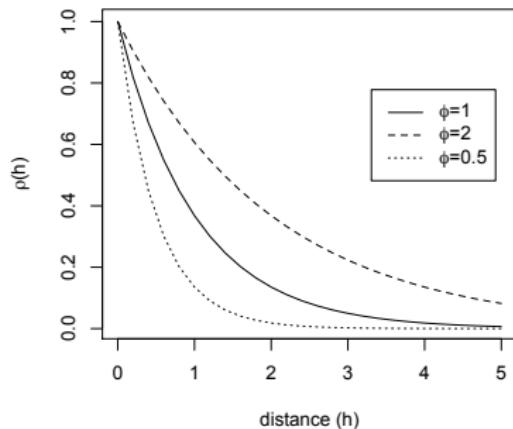
Everything is related to everything else, but near things are more related than distant things (first law of geography)

Waldo Tobler, 1970

Moreover . . . typically, beyond a certain distance between two observations, they become (spatially) uncorrelated/independent.

How to Control the Spatial Correlation?

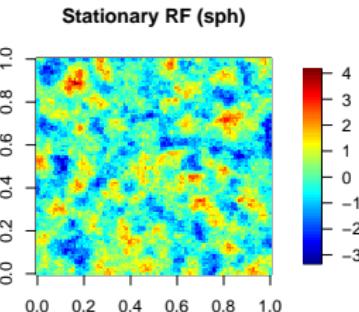
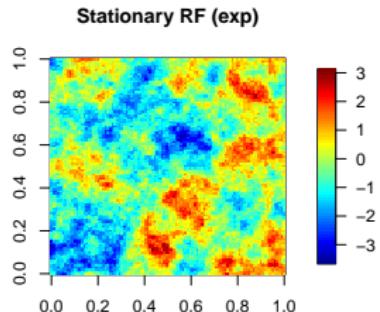
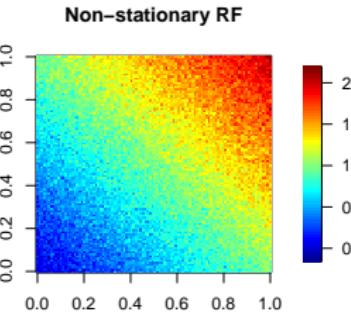
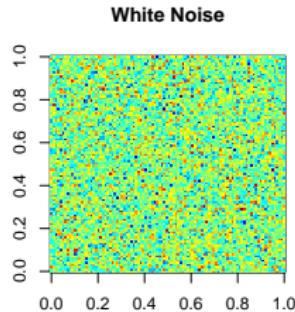
Under certain stationarity (stability) conditions:



Left Panel: Exponential correlation function for different values of the [radius of influence](#), also called *range*, ϕ .

Right Panel: Schematic representation of a [variogram](#) with its structural parameters.

Simulations of Geostatistical Data (also named Random Fields, RF)



Outline

① Introduction to Spatial Statistics

- ▷ Data Analysis and Visualization

② Spatial Data: **Areal, Point, Geostatistical**

- ▷ Understanding Spatial Correlation

③ Applications to Fisheries and Marine Sciences

- ▷ Try DYNAMIC MAPS

Marine Spatial Planning (MSP)

- MSP started in the 1980's, with involvement of ICES².
- Conflicts arise in **resource usage, environmental impact, and legal frameworks**.



Figure: A crowded ocean under a changing climate. Cartoon created by Bas Kohler (Santos et al., 2020).

²International Council for the Exploration of the Sea (ICES)

The Portuguese Case – IPMA / ICES

IPMA is a public institution under the Ministry of Economy and Sea, supporting the **strategic plans in term of analysis and survey of Marine Resources**.

Key Objectives:

- ① Combating marine biodiversity loss, whose main causes are
 - Climate changes
 - Illegal fishing activities
 - Difficult management of fishery resources
- ② Improving knowledge on the population dynamics of marine species with commercial interest, such as sardines, anchovies, hake, etc.



recovery of marine biodiversity
+
sustainability of fisheries

How can Spatial Statistics help?

Creating mathematical/statistical models for . . .

- **Space-time analysis** of abundance indicators. Examples:
 - What are the best places to fish for sardines? And at what time of the year?
- **Studying the influence** of environmental/biological variables. Examples:
 - Do sardines prefer salty water? And what temperature?
- **Construction of maps.** Examples:
 - What are the preferred locations for juvenile individuals?

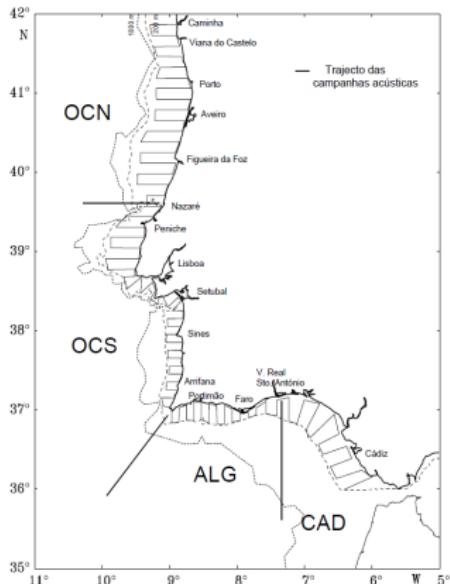
Summarizing. . . spatial statistics helps to assess the status of fish populations and make informed management decisions!

Data sources

Scientific surveys



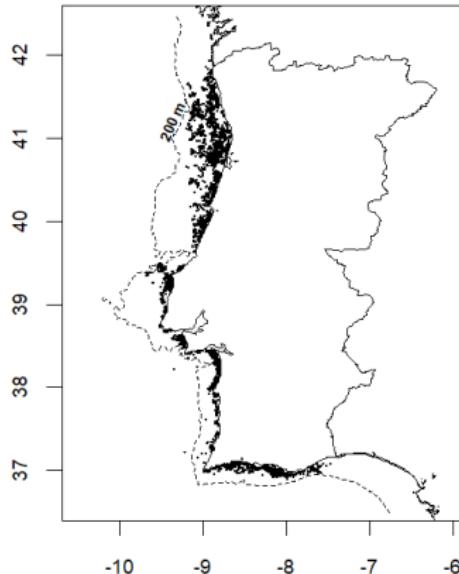
Fishery-independent data



Commercial fisheries

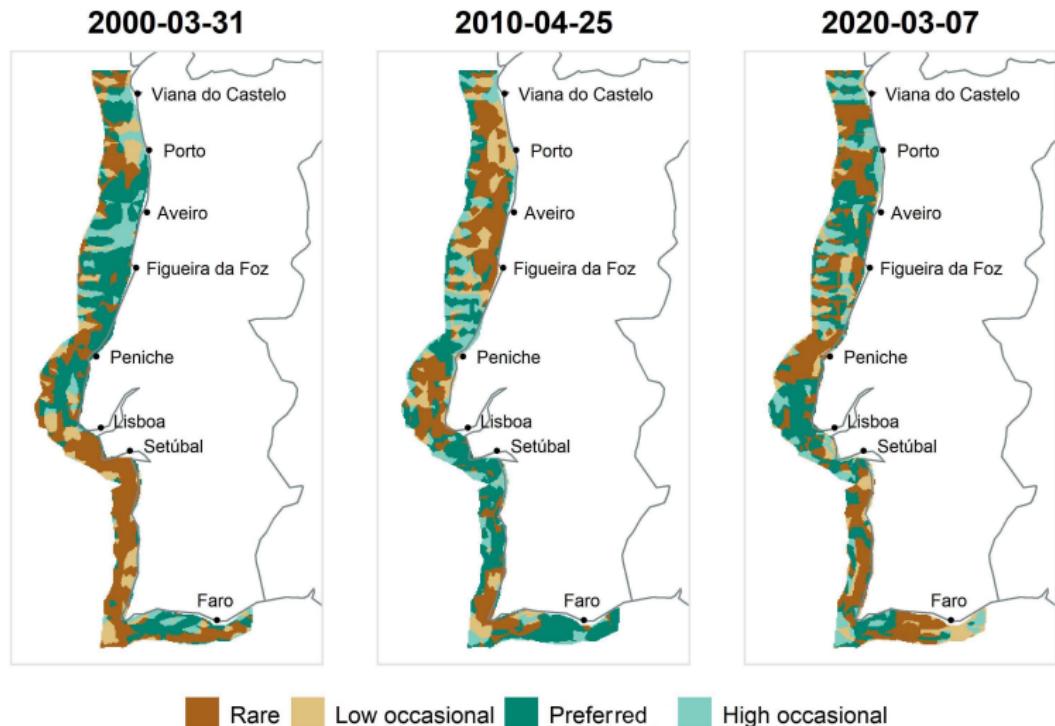


Fishery-dependent data



Spatial Statistics Supporting Decision Making

Areas of occupancy for sardine



Thank you for your time and interest!

Feel free to ask any questions – and explore the dynamic maps at your own pace after the talk.

Instructions:

- ① Use your phone to scan the QR code.
- ② Open the HTML file found on Raquel's GitHub.
- ③ Explore the interactive boxplots and other graphics.
- ④ **Navigate and zoom in on the world map to find the coast of Portugal... and data from the scientific campaigns.**

