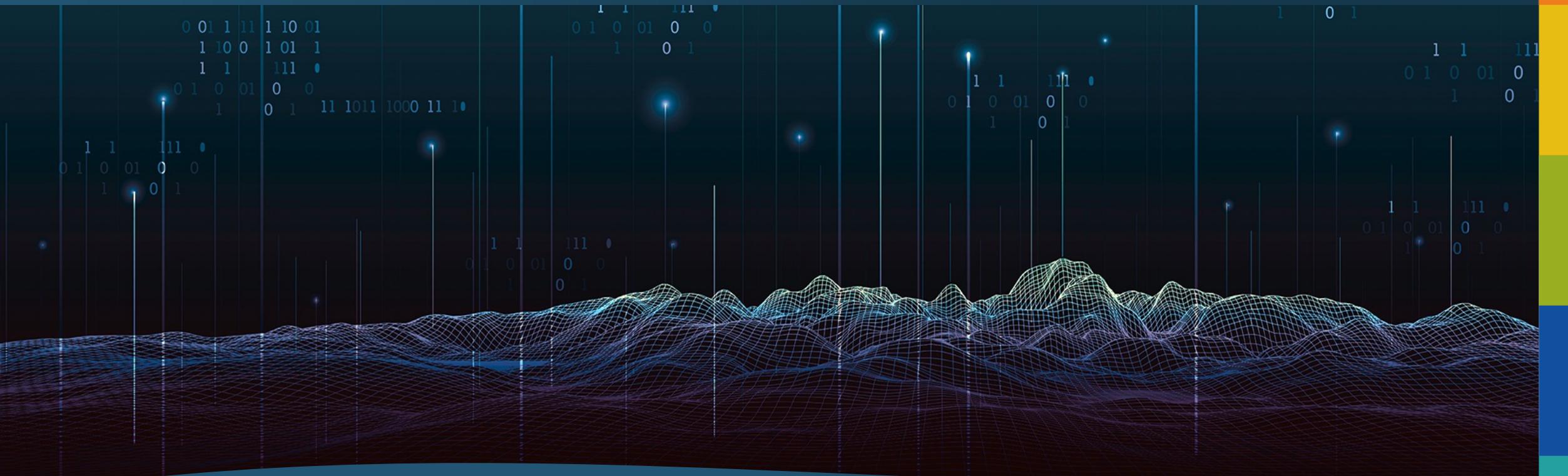


Spatial data and data integration

Australian Urban Research Infrastructure Network (AURIN)



Acknowledgement of Country



We acknowledge the Traditional Owners of the land on which this event is taking place and pay respect to their Elders (past and present) and families.

The Geosocial work package

Finding and using spatial data

Deciding on integration

Integration methods

Producing a new data product

Exercise: Producing a map

The Geosocial work package

Finding and using spatial data

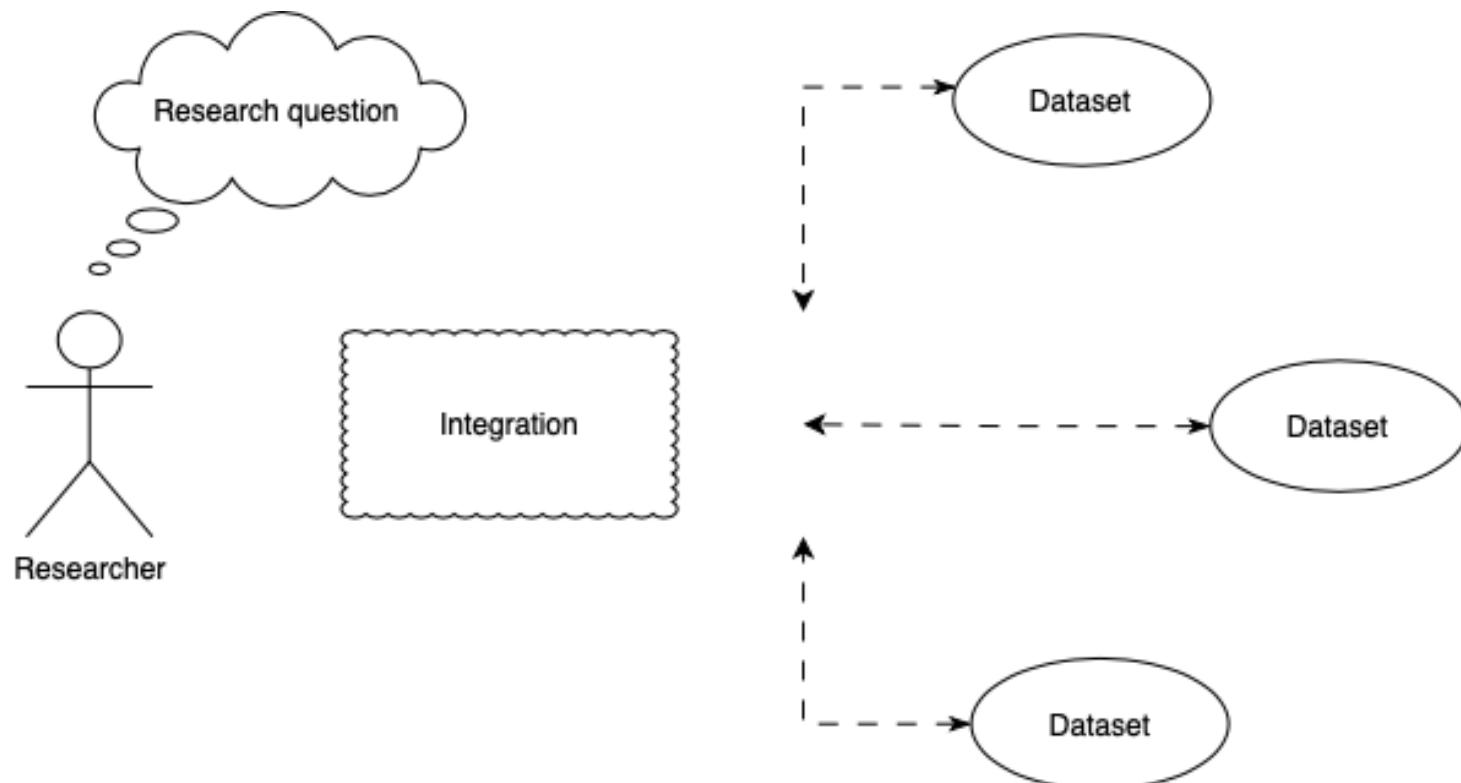
Deciding on integration

Integration methods

Producing a new data product

Exercise: Producing a map

Problem: The researchers want to bring data on people and places together, but don't know how to do it and what the issues might be.



The Geosocial work package: Motivation



Solution: Data integration service which will allow researchers to enhance **people-centred survey data with spatially structured data** capturing information on places where these people live.



The Geosocial work package

Finding and using spatial data

Deciding on integration

Integration methods

Producing a new data product

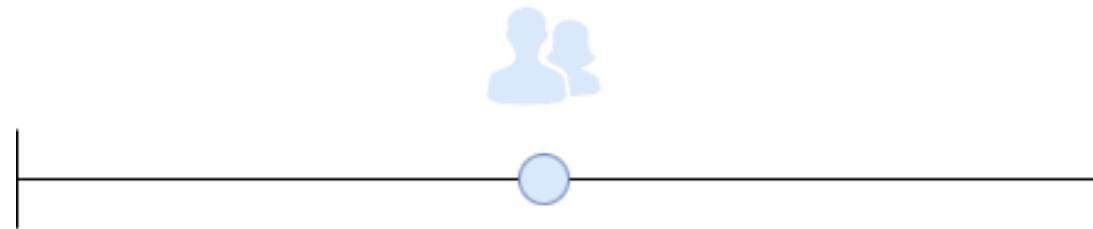
Exercise: Producing a map

Survey data

Cross-sectional data: Cross-sectional data consists of data on one or more variables collected at the same point in time. (Gujarati, 2011)

Examples:

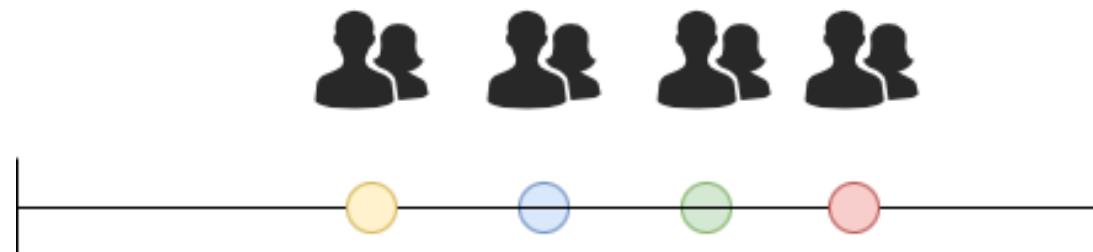
- Population census
- Consumer Expenditure Surveys
- Opinion polls



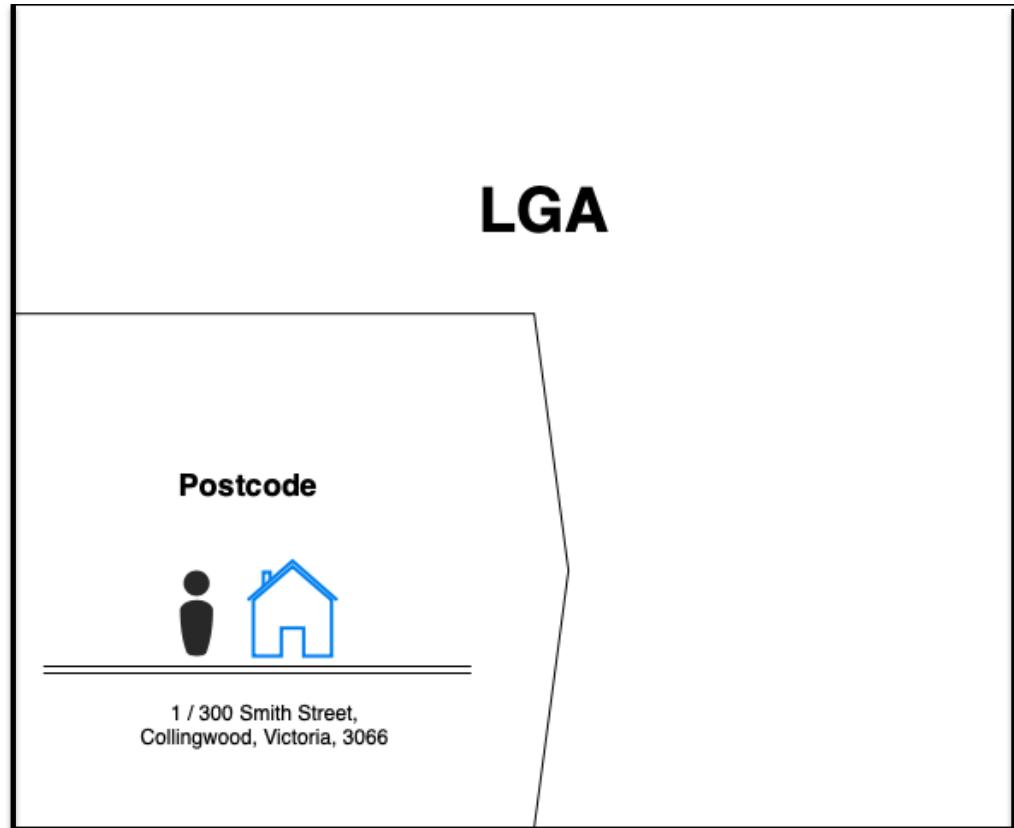
Longitudinal data (panel data): repeatedly collect data from the same sample over an extended period of time

Examples:

- Household, Income and Labour Dynamics in Australia (HILDA)
- Longitudinal Surveys of Australian Youth (LSAY)

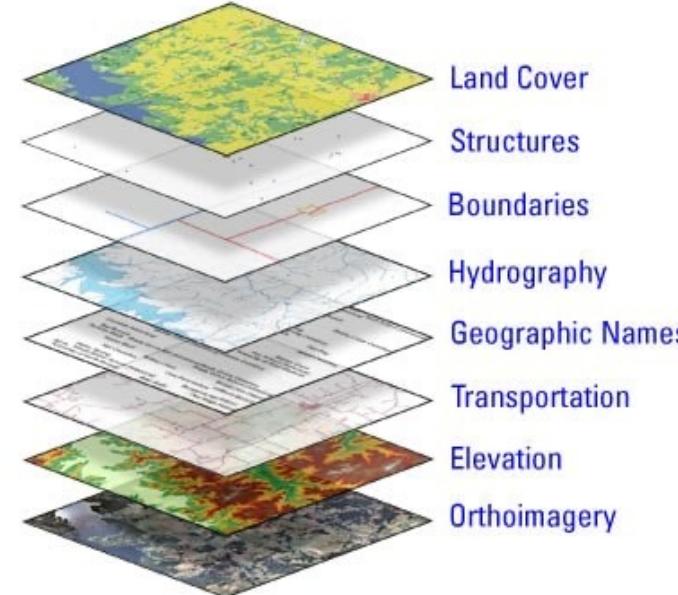


What is spatial data?

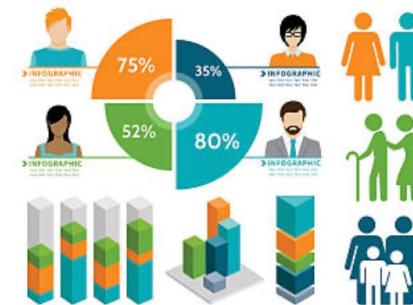


Data that have an implicit or explicit association
with a location relative to Earth

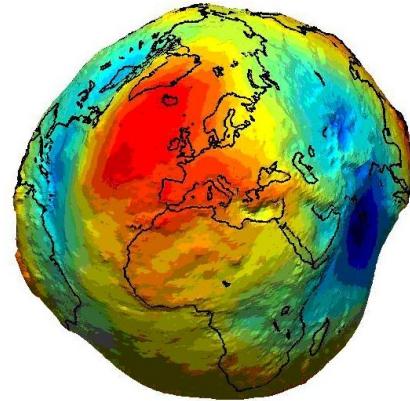
Physical characteristics



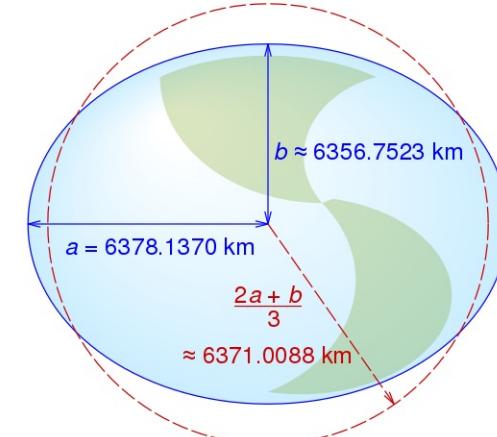
Socio-demographic characteristics



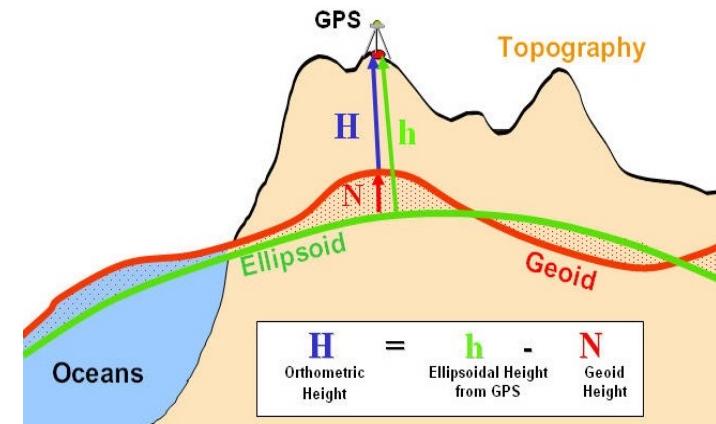
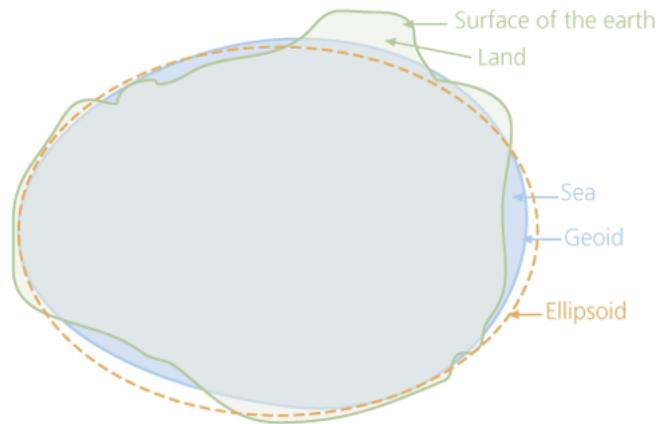
World Geodetic System 1984 (WGS84)



Geoid



Ellipsoid



Geoid vs Ellipsoid

Finding and using spatial data

Spatial analysis: represents a collection of **techniques** and **models** that explicitly use the spatial referencing of each data case.

Spatial analysis needs to make assumptions about or draw on data describing spatial relationships or spatial interactions between cases. (Chorley, 1972; Haining 1994).



The Geosocial work package

Finding and using spatial data

Deciding on integration

Integration methods

Producing a new data product

Exercise: Producing a map

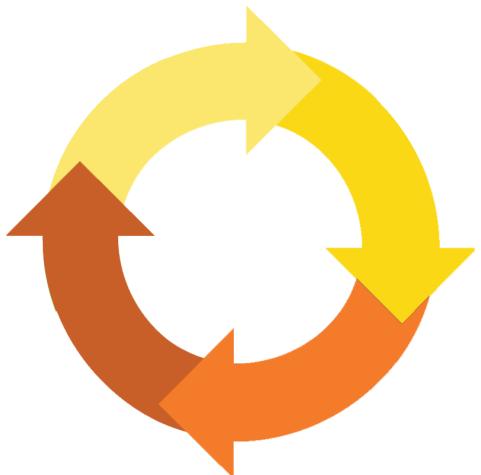
Deciding on integration:

Step 1: Decide a research question.

Step 2: Find datasets that are suitable to the research question.

Step 3: Define a methodology to answer your research question

Step 4: Robustness: Evaluate the methodology



Deciding on integration:

Step 1: Research question: How is leaving a high school associated with Australia's economic development?

Step 2:

Database: PHIDU - Education (LGA) 2015-2016:

Variable: School Leaver Participation In Higher Education 2016 Enrolled in higher education.

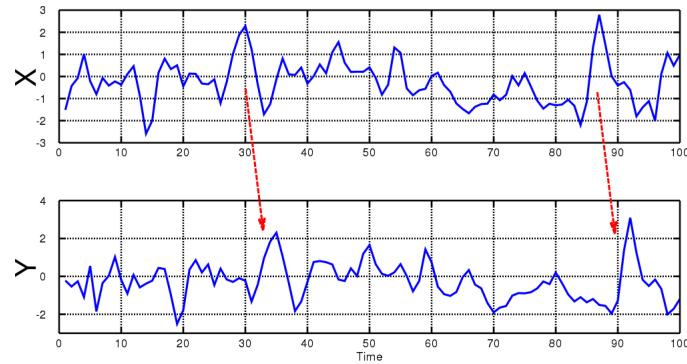
Database: Socio-Economic Indexes for Areas (SEIFA) – SA2 - 2016

Variable:

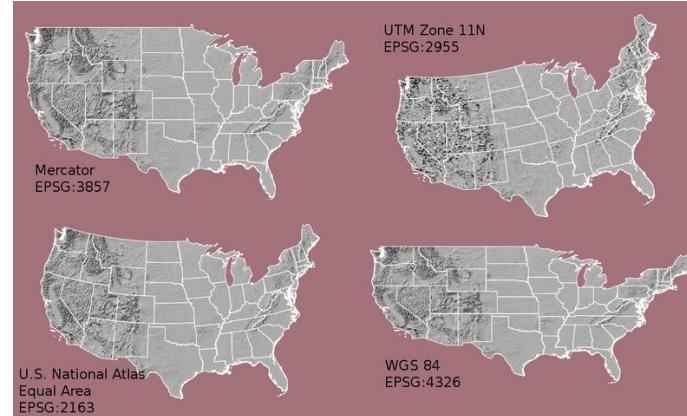
Step 3: Use a spatial correlation (Moran's I) to evaluate the correlation between leaving the school on economic development (SEIFA).

Step 4: Try other ways to estimate spatial correlation: Geary's C and Getis-Ord's index

Considerations



Causality/ Temporary lags

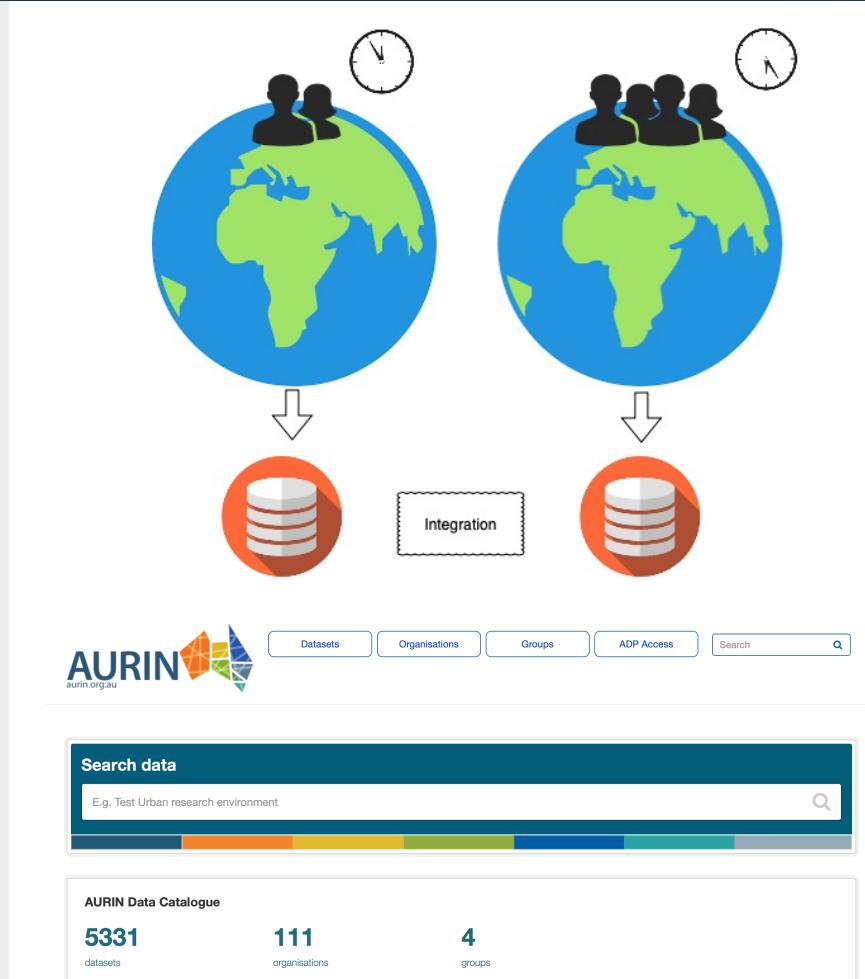


Spatial



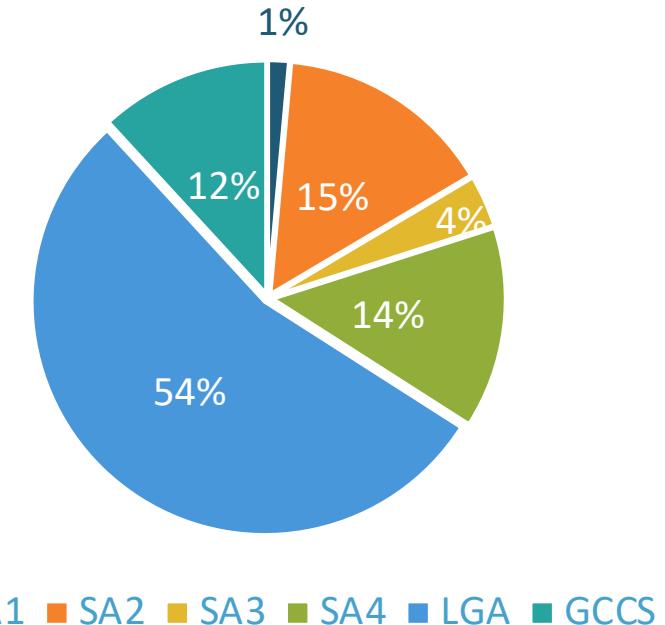
Semantic

Deciding on integration:



social geography youth health
human society
safety polygon
geography human
vic education policies human
society human health development
education children work
society demography
victoria vic
society geography economic
economic employment demographics socio
population trends
economic geography
polygon victoria
human geography
children youth geography social
demography population
economics economic
employment demographics development safety
socio economics

■ SA1 ■ SA2 ■ SA3 ■ SA4 ■ LGA ■ GCCS



More information: <https://data.aurin.org.au>

The Geosocial work package

Finding and using spatial data

Deciding on integration

Integration methods

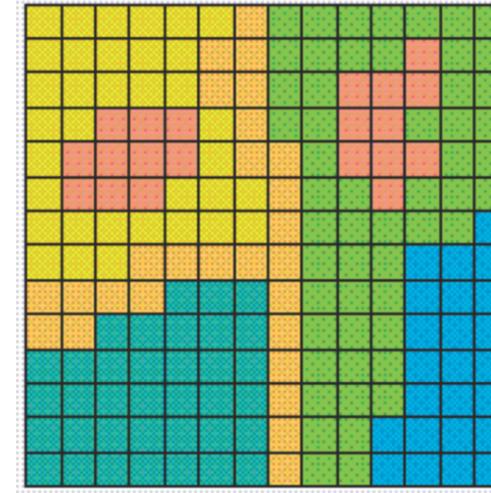
Producing a new data product

Exercise: Producing a map

Finding and using spatial data

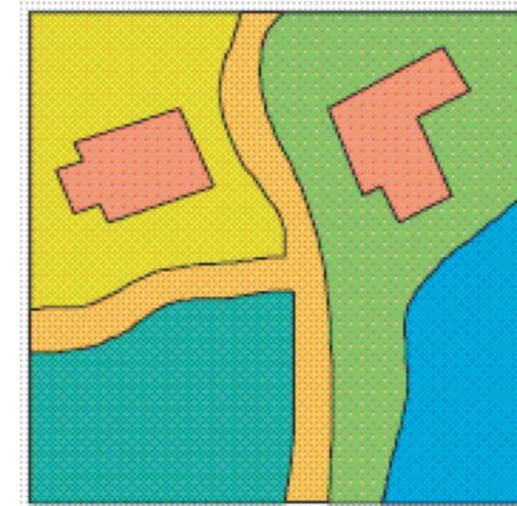
Raster map:

- Raster data is stored as a grid of values that are rendered on a map as pixels.
- Each pixel value represents an area on the Earth's surface.



Vector map:

- Consists of objects described by coordinates in a given coordinate system.
- The vector model uses points and line segments to identify locations on the earth.



Vectorial map:

Type of elements:

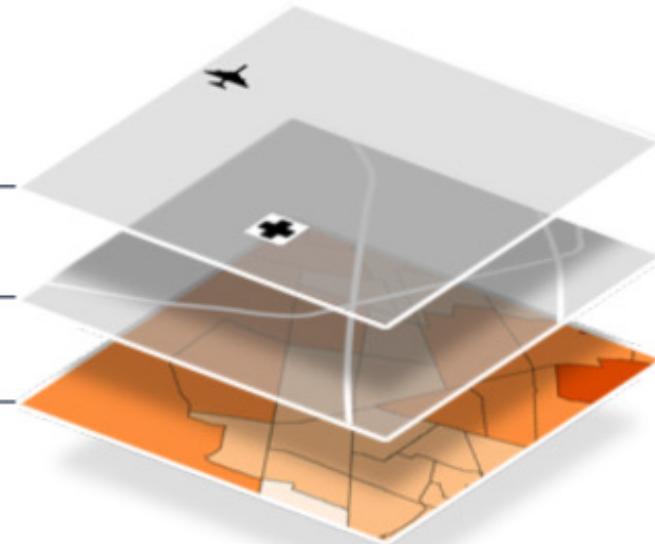
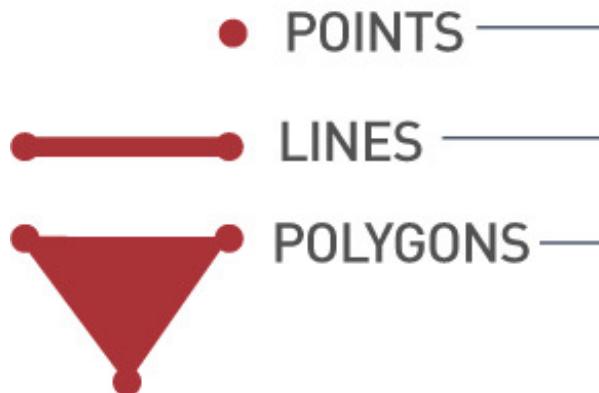
Point: Addresses, locations, points of interest, etc

Lines: streets, freeways, borders, etc

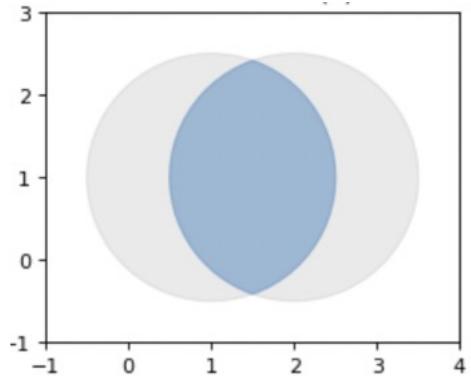
Polygons: Countries, cities, Cadastre

Advantages:

- Spatial operations
- Spatial aggregation
- Spatial Join



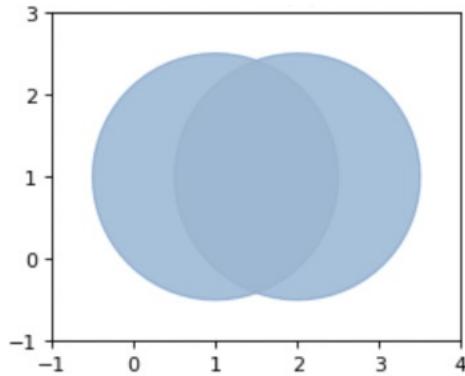
Spatial operations



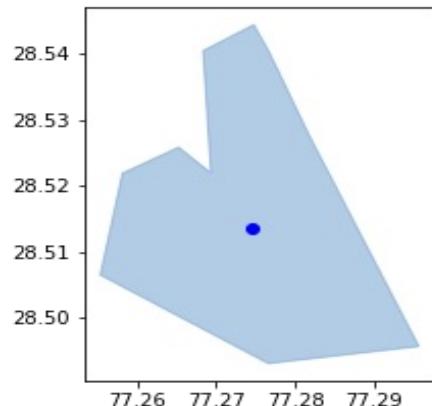
Intersection



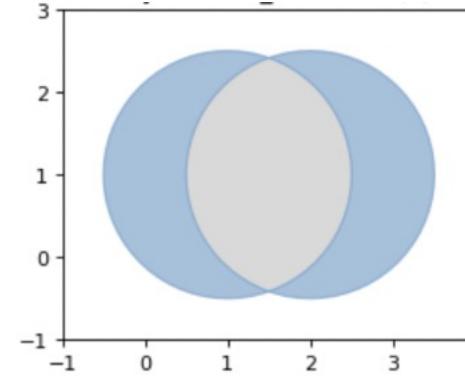
Contour



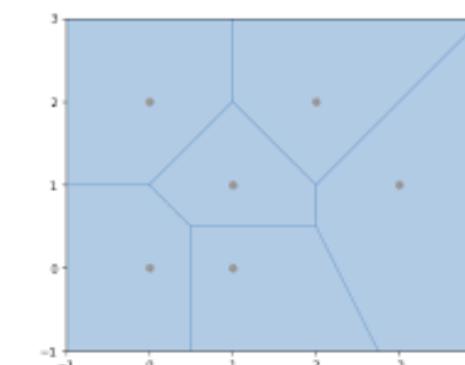
Union



Difference

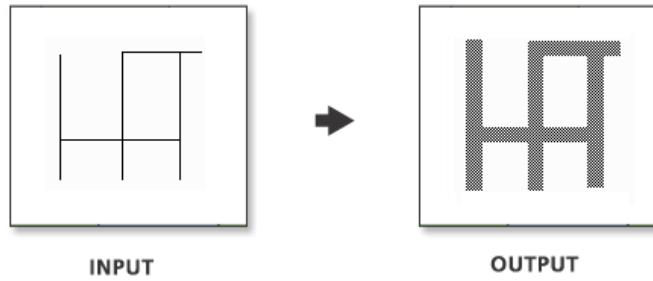


Centroid

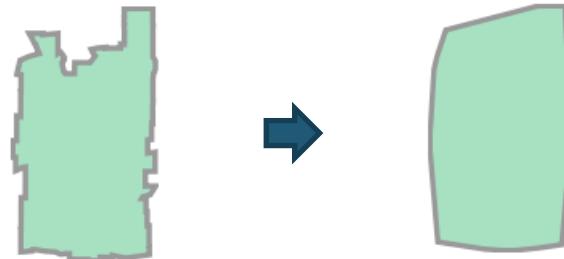


Images from: Pysal

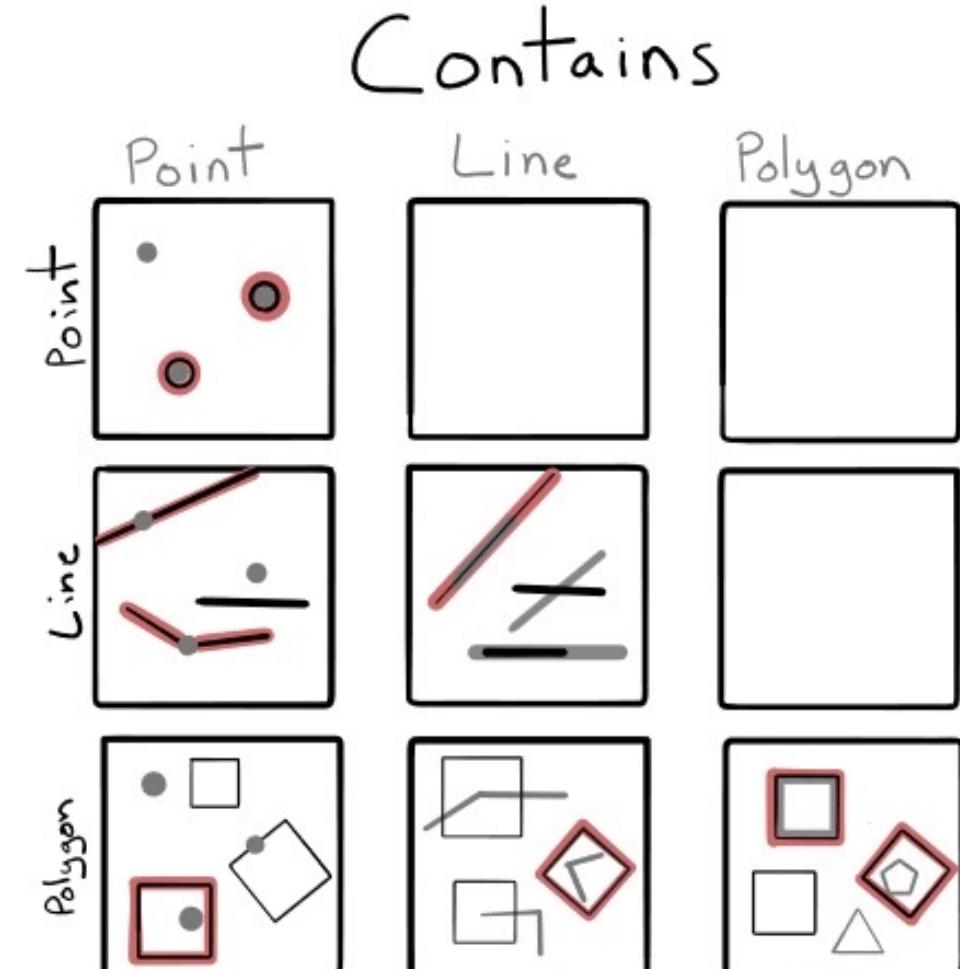
Spatial operations



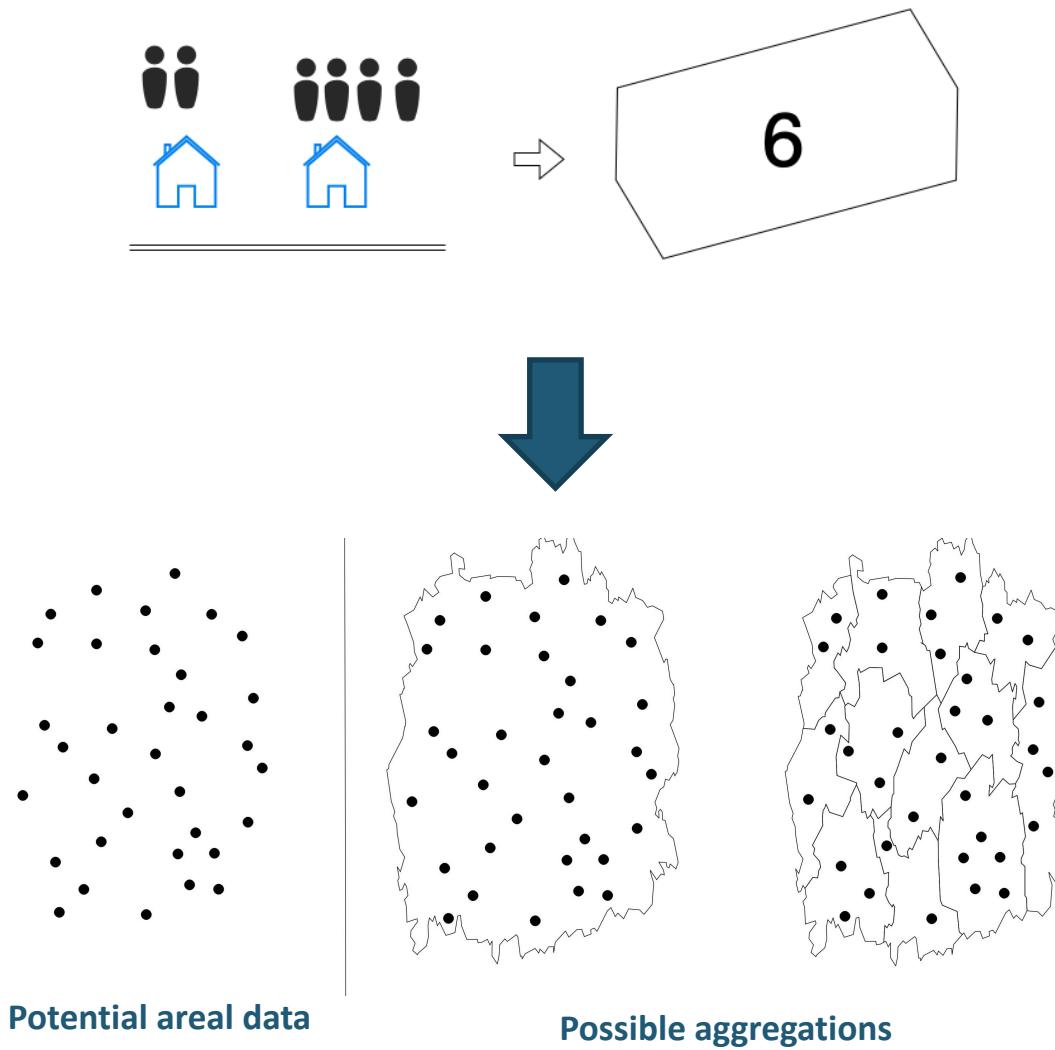
Buffer: Create a new layer that covers this in a zone of influence whose radius is the one indicated in the analysis tool.



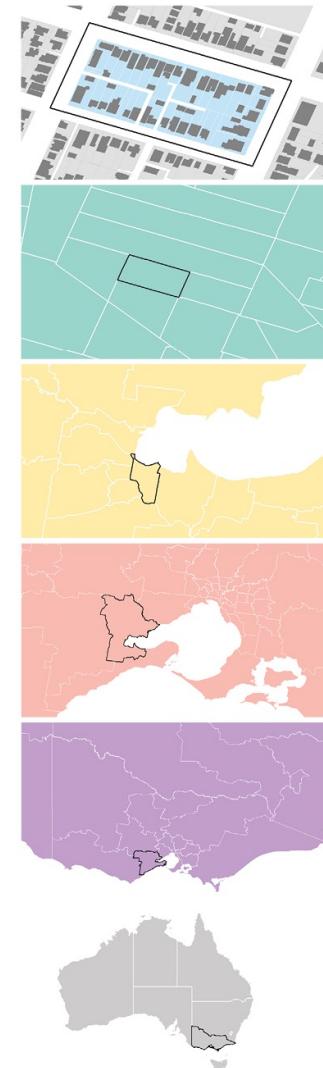
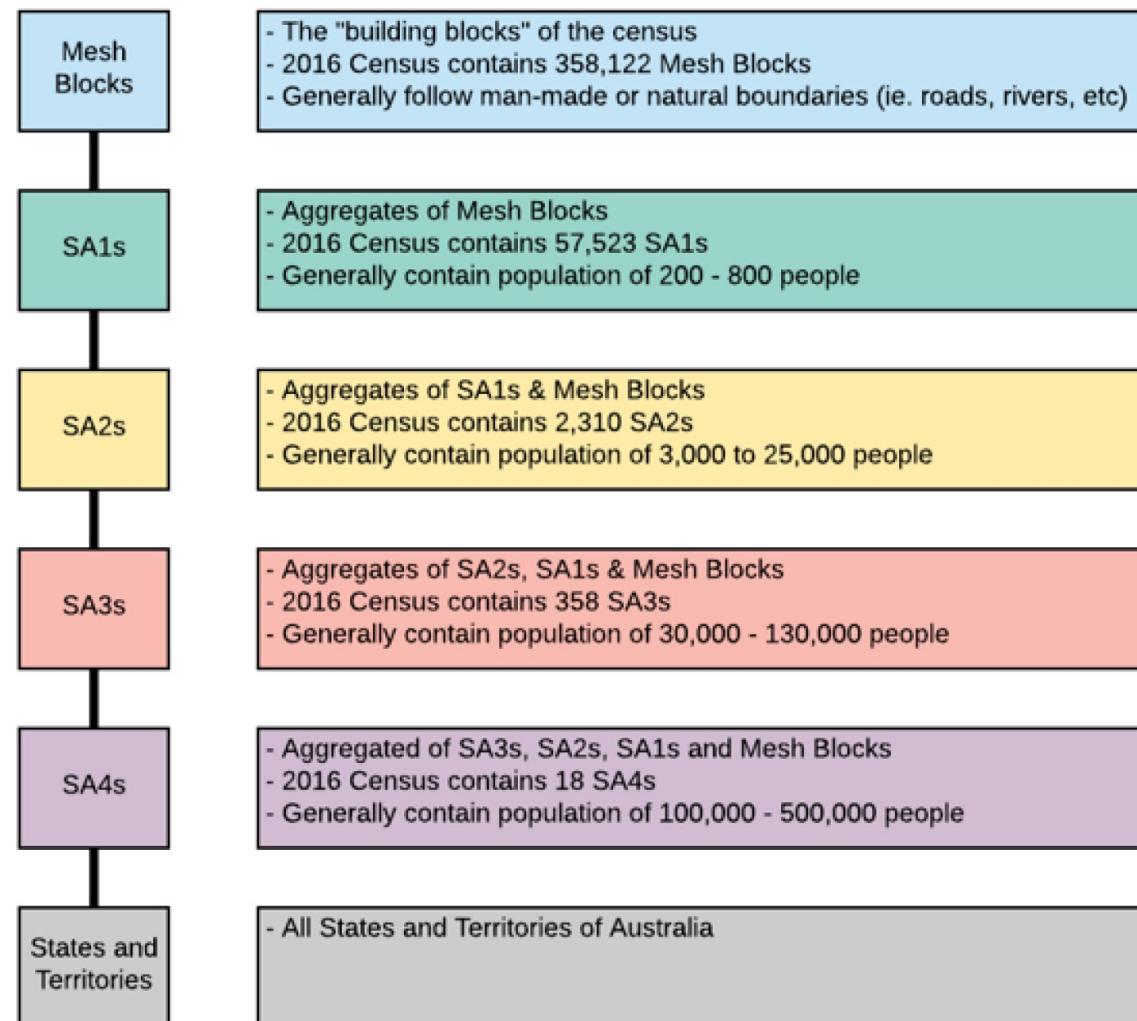
Convex capsule: a set of points that contains the intersection of all convex sets within a polygon.



Spatial aggregation



Spatial aggregation



Concordance/Correspondence:

Objective: mathematically convert statistical data to and from geographic regions

Main Structure and Greater Capital City Statistical Areas

2016 Mesh Blocks to 2021 Mesh Blocks

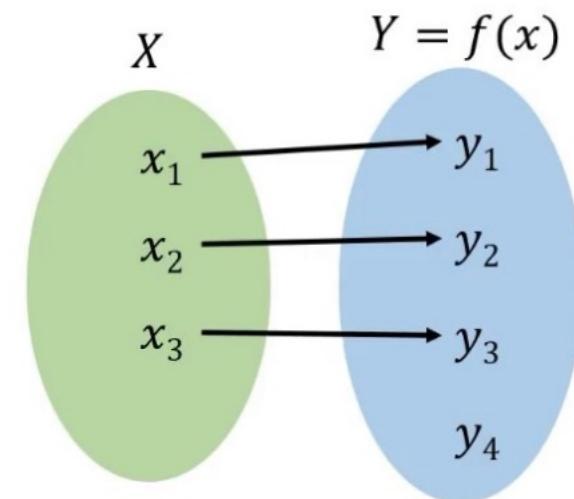
2016 Statistical Areas Level 1 to 2021 Statistical Areas Level 1

2016 Statistical Areas Level 2 to 2021 Statistical Areas Level 2

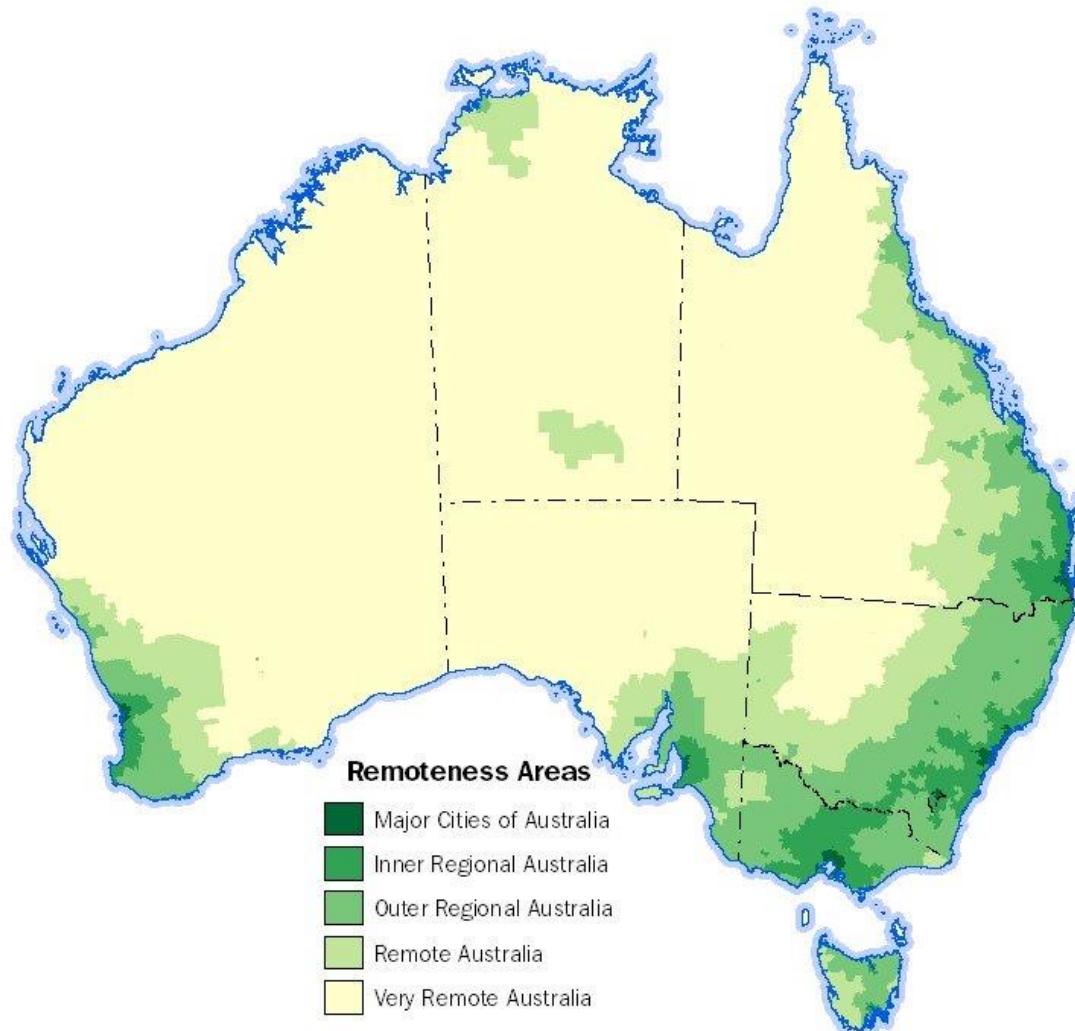
2016 Statistical Areas Level 3 to 2021 Statistical Areas Level 3

2016 Statistical Areas Level 4 to 2021 Statistical Areas Level 4

2016 Greater Capital City Statistical Areas to 2021 Greater Capital City Statistical Areas



Spatial aggregation: Choropleth map

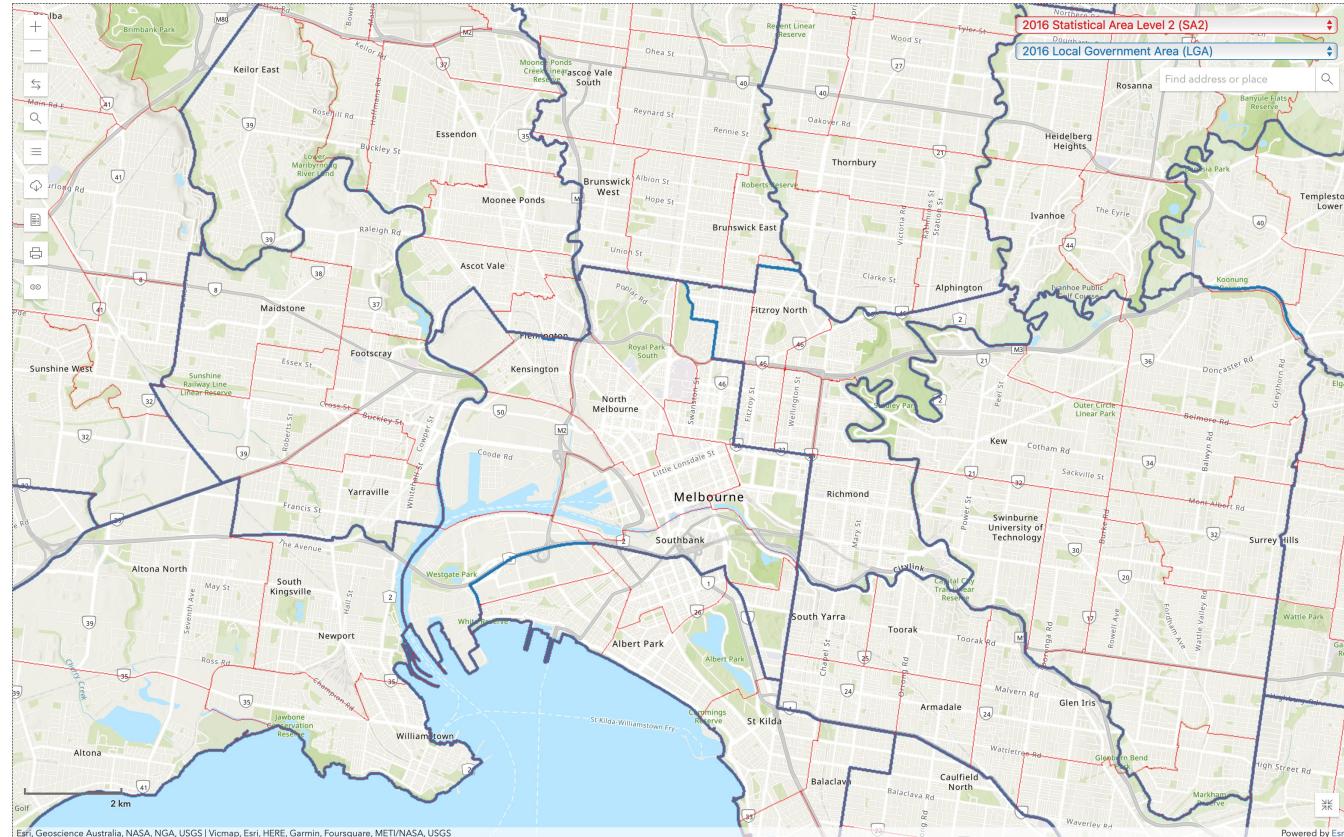


Example: Spatial aggregation

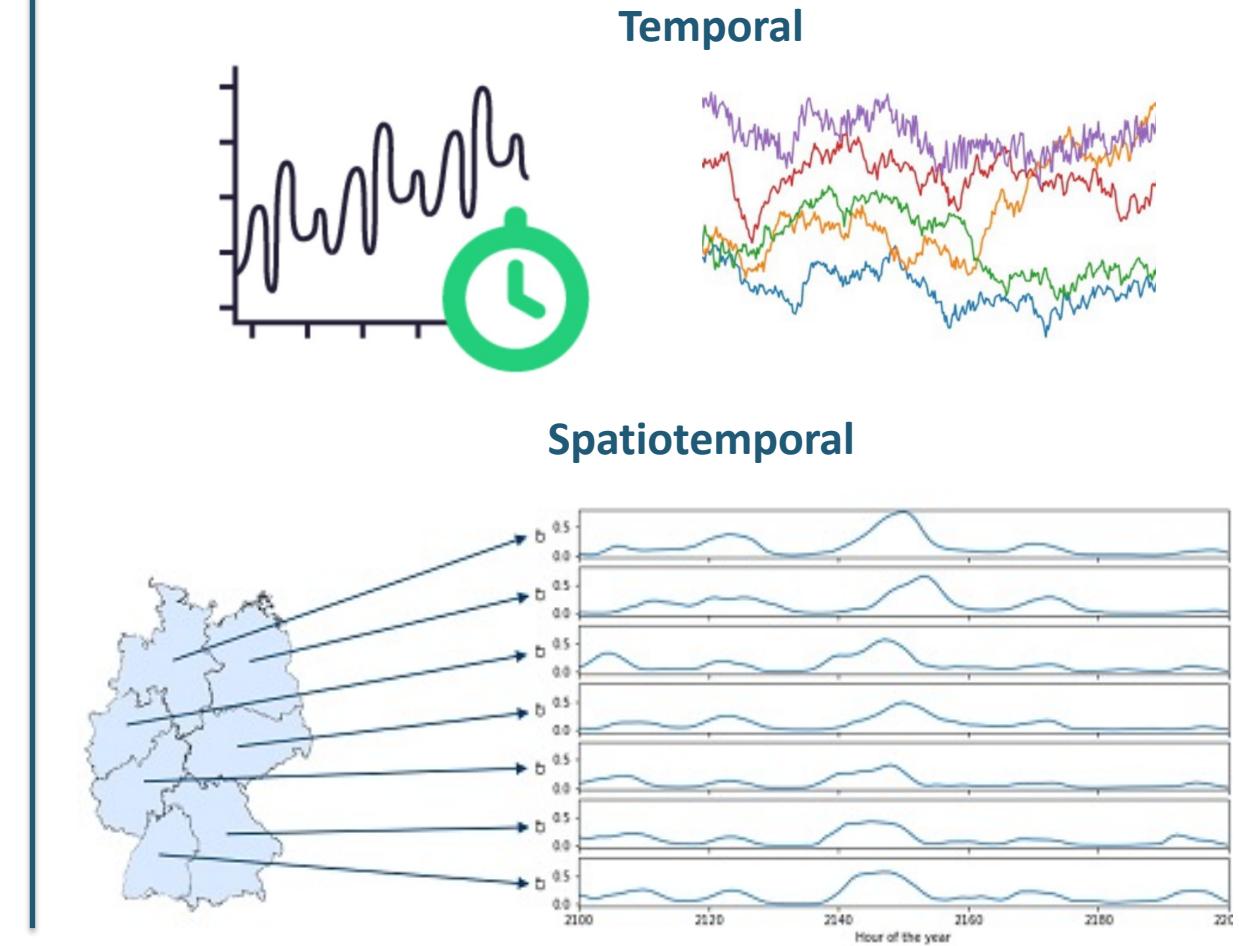
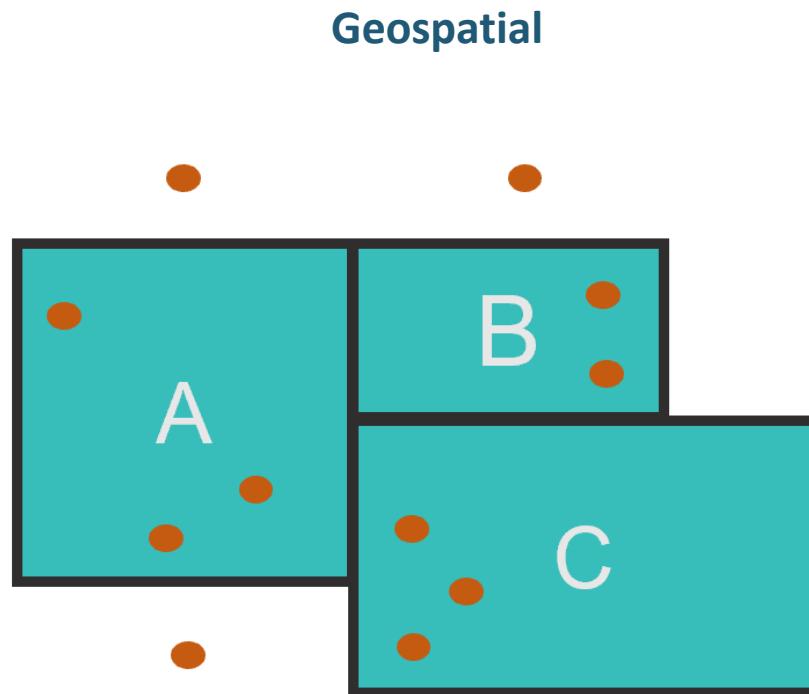
Socio-Economic Indexes for Areas (SEIFA) – SA2 - 2016



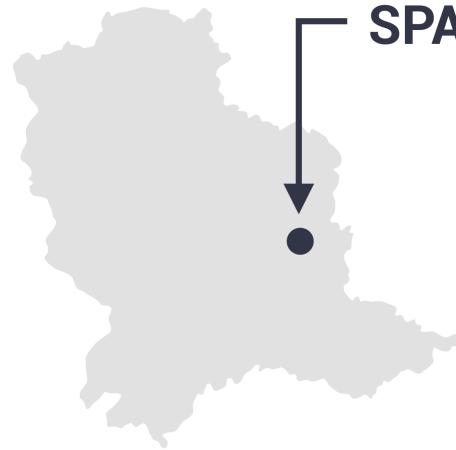
PHIDU - Education (LGA) 2015-2016:



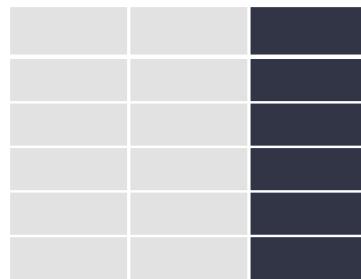
Spatial Join



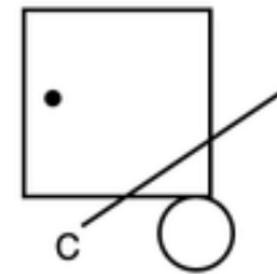
Spatial Join



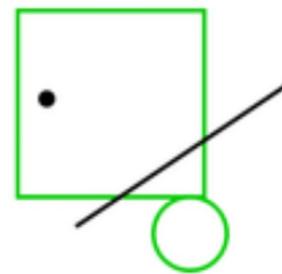
SPATIAL JOIN



Types of spatial join

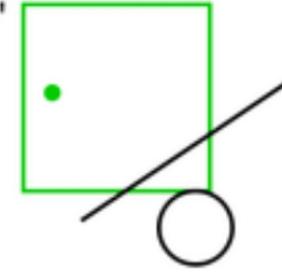
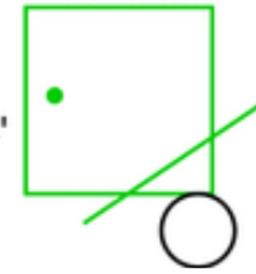


spacial join:
op = 'touches'



spacial join:
op = 'contains'

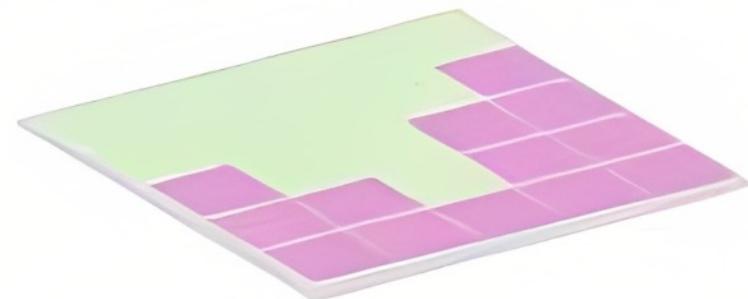
spacial join:
op = 'intersects'



Spatial autocorrelation - Example

"Everything is related to everything else, but near things are more related than distant things."

Waldo Tobler



High Spatial Autocorrelation
(Clustering)



Low Spatial Autocorrelation
(Checkerboard)

Limitations

- **Spatial aggregation:** In some cases, it is impossible to delve into the smallest detail of the problem.
- **Measurement error:** Data that does not correspond to the true values.
- **Assumptions:** Assumptions that are not entirely realistic.
- **Computing capacity:** High computational costs.



The Geosocial work package

Finding and using spatial data

Deciding on integration

Integration methods

Producing a new data product

Exercise: Producing a map

Producing a new data product

A report, dashboard or application with its own user interface (UI), an API, or command-line SQL access



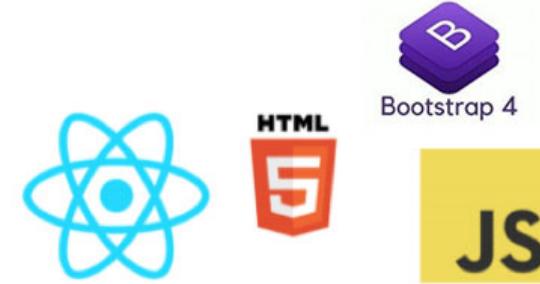
Image from: "Product Data and Your Digital Catalog" by Handshake

Producing a new data product



```
def add5(x):
    return x+5

def dotwrite(ast):
    nodename = getNodename()
    label=symbol.sym_name.get(int(ast[0]),ast[0])
    print '%s %s %s' % (nodename, label),
    if isinstance(ast[1], str):
        if ast[1].strip():
            print '%s';% ast[1]
        else:
            print ''
    else:
        print ']';
        children = []
        for n, child in enumerate(ast[1:]):
            children.append(dotwrite(child))
        print '%s -> [%s' % nodename,
        for name in children:
            print '%s' % name,
```



```
1 <!DOCTYPE html>
2 <html>
3     <head>
4         <title>Example</title>
5         <link rel="stylesheet" href="styl
6     </head>
7     <body>
8         <h1>
9             <a href="/">Header</a>
10        </h1>
11        <nav>
12            <a href="one/">One</a>
13            <a href="two/">Two</a>
14            <a href="three/">Three</a>
15        </nav>
```

Examples:

- R - Shiny: <https://shiny.rstudio.com/gallery/>
- Python Dash: <https://dash.gallery/Portal/>
- Java - <https://kepler.gl>

The Geosocial work package

Finding and using spatial data

Deciding on integration

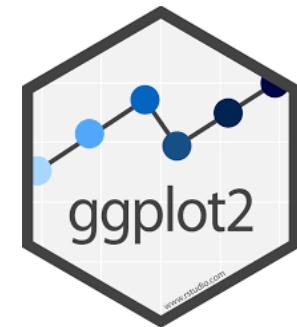
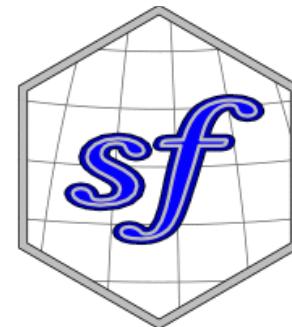
Integration methods

Producing a new data product

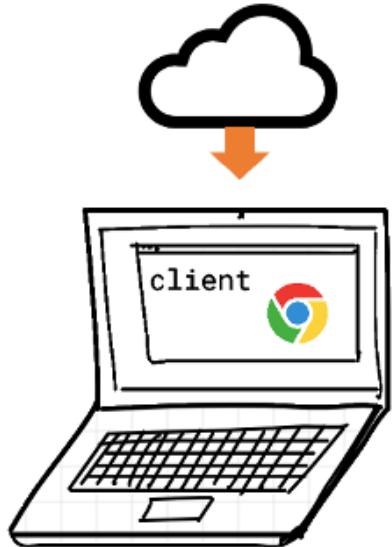
Exercise: Producing a map

Exercise – R

- **Part 1:** Spatial operations
 - Union
 - Intersects
 - Difference
 - Within
 - Buffer
 - Convex hull
 - Centroid
 - Boundaries
 - Area (Mts²)
- **Part 2:** Integration methods:
 - ID join: Education & LGA by LGA_ID
 - Spatial join: SEIFA – LGA database
- **Part 3:** Geospatial aggregation:
 - SEIFA → SA2 → LGA
- **Part 4:** Producing a map

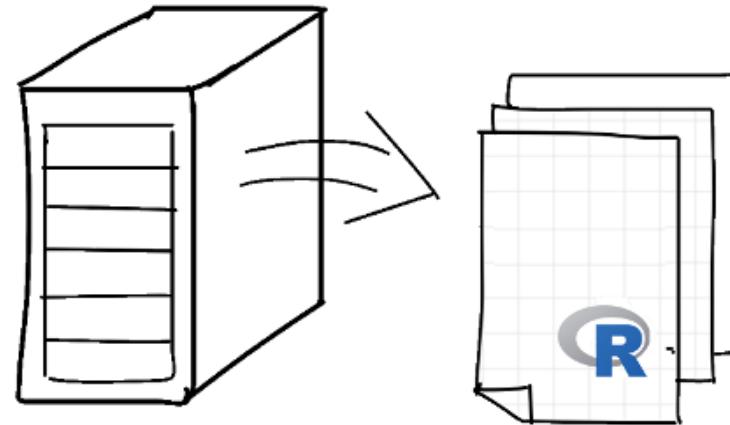


Choices for the R environment



a) Cloud environment

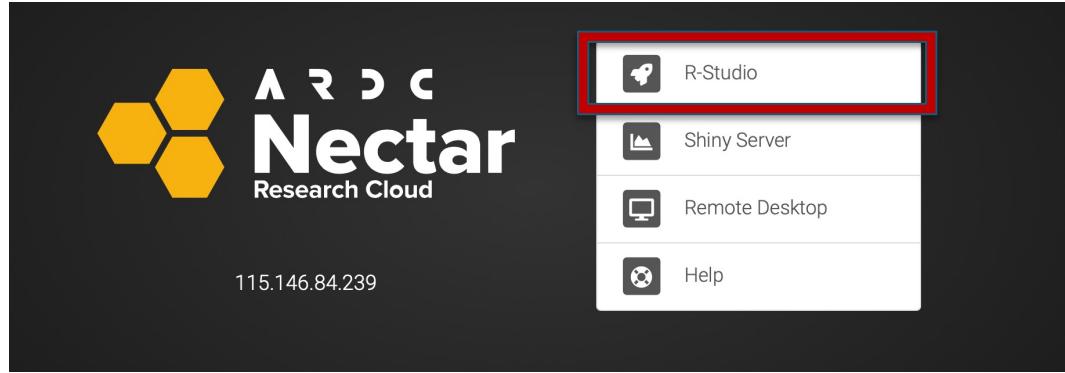
- Compatibility
- Less time spent installing R and dependencies.
- Better Package compatibility
- New users



b) Local environment

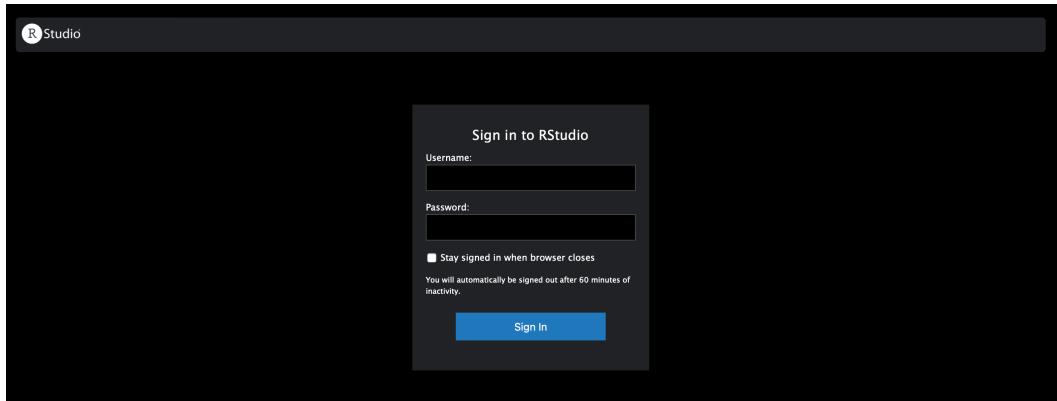
- Customisable
- Privacy
- Large data sets
- Experienced users

Cloud environment



Group 1:

- Host: <http://203.101.230.180>
- User: user1....user15
- Password: hass2023



Group 2:

- Host: <http://203.101.231.227>
- User: user1....user15
- Password: hass2023



Local environment



Step 1: Clone the repository on GitHub

https://github.com/AURIN-OFFICE/HASS_Summer_School

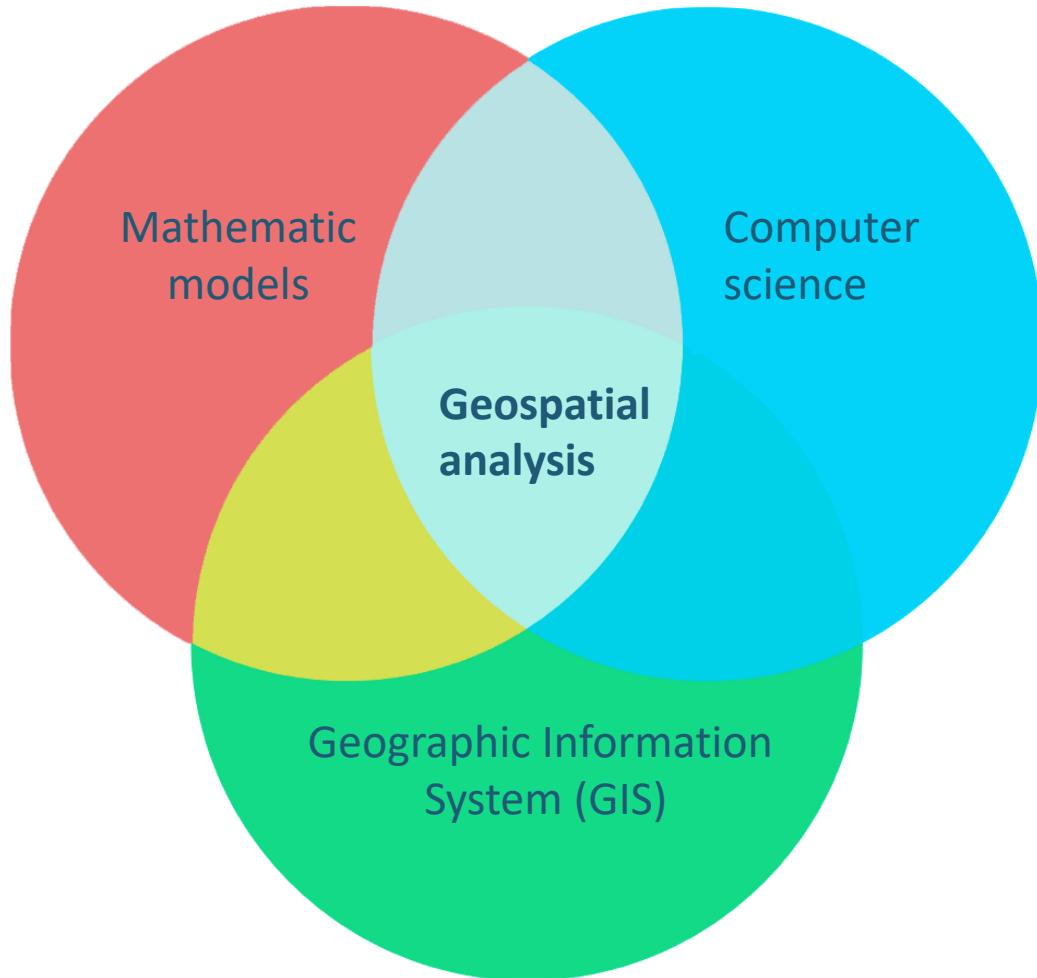


Step 2: Install the necessary libraries

```
##### ----- Workshop ----- #####
##### ----- Clean variables ----- #####
rm(list=ls())
##### ----- Install libraries ----- #####
# install.packages(c('sf', 'tidyverse', 'ggplot2', 'leaflet.extras'))
```

Thank you

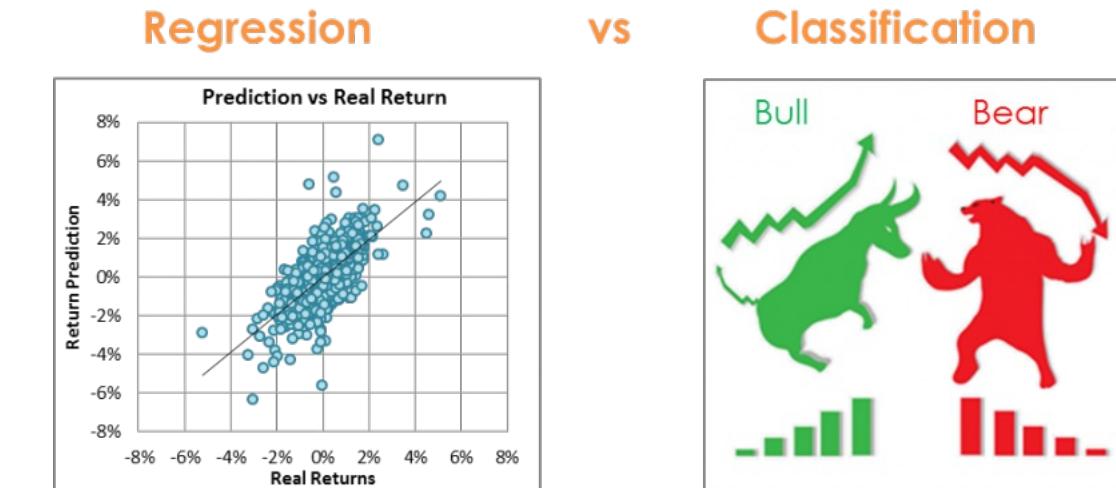
Finding and using spatial data



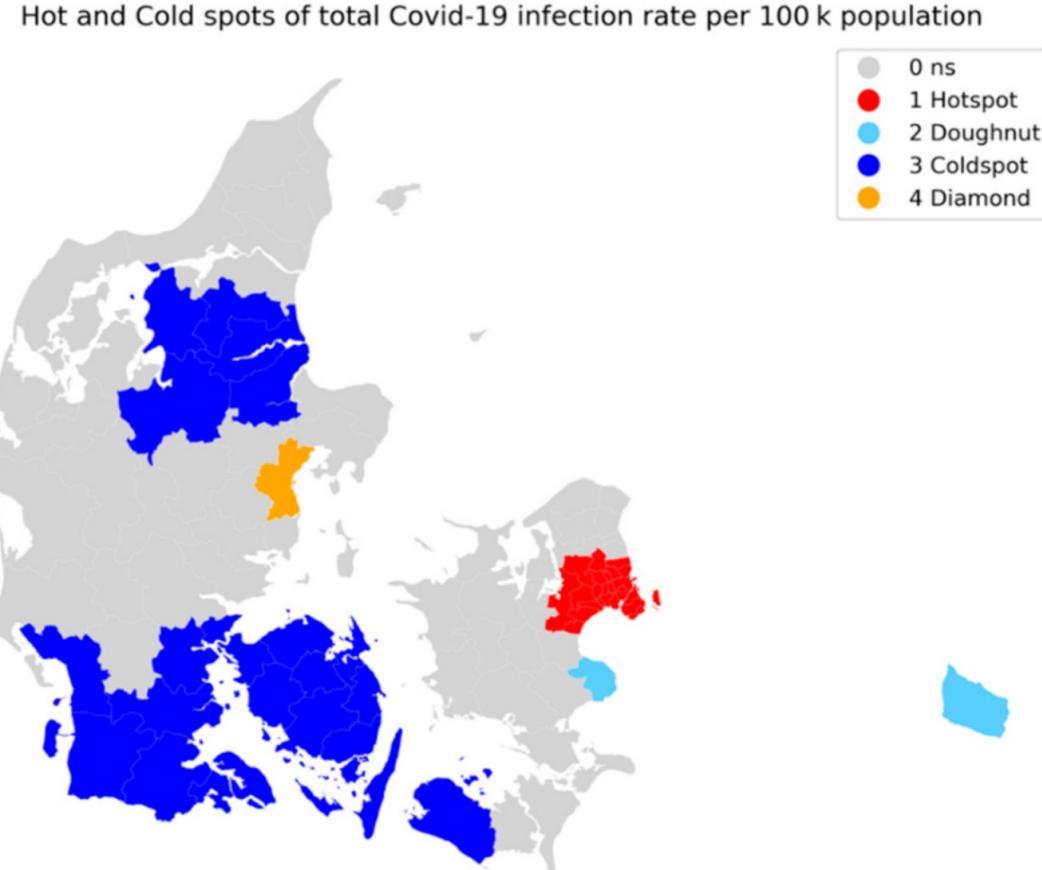
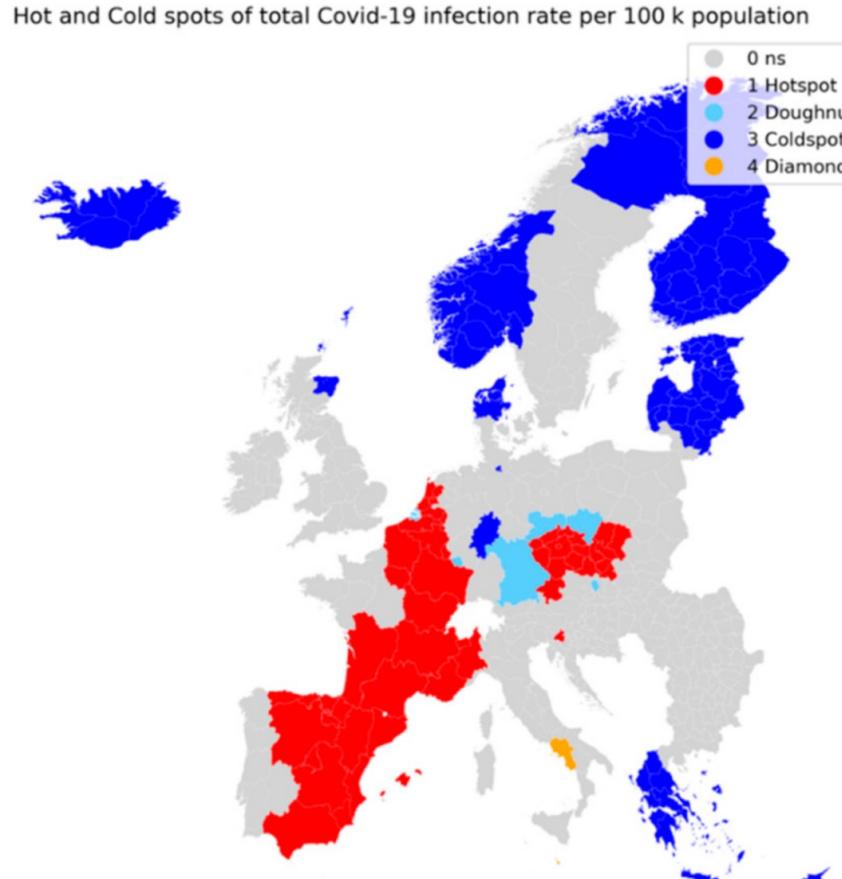
Supervised vs unsupervised learning

Supervised models:

Supervised learning uses as training a data set that contains a mark or label on the data. In this type of learning, the distinction between the independent and dependent variables is clear.



Spatial autocorrelation - Example



Source: Hass, Frederik Seeup, and Jamal Jokar Arsanjani. "The geography of the COVID-19 pandemic: A data-driven approach to exploring geographical driving forces." International Journal of Environmental Research and Public Health 18.6 (2021): 2803.

Unsupervised models

Unsupervised models

Unsupervised learning does not require prior labelling of the data, it uses all the information to make associations between the data or group them.

