UNIVERSITY OF TWENTE.

Introduction to Exploratory (Spatial) Data Analysis

Mahdi KHODADADZADEH Assistant Professor Faculty of Geo-Information Science and Earth Observation (ITC) Department of Geo-information Processing (GIP) m.khodadadzadeh@utwente.nl

May 2024



FACULTY OF GEO-INFORMATION SCIENCE AND EARTH OBSERVATION





From: https://xkcd.com

This lesson's learning objectives

Explain to peers

- the fundamentals of E(S)DA
- the importance of E(S)DA before modelling

Apply statistical and visualization methods on different types of data

Develop familiarity with Python





You are a Python master. Congrats!





You've learned how to build a model in Python. Congrats!





But you run ínto some íssues!



Data Analysis Workflow



From: https://davpy.netlify.app/3-data-workflow.html



Ingesting Data

Getting data in a shape that we can use to start our analysis.

- Python:
 - Reading comma separated value (CSV) data: pandas.read_csv()
 - Reading an Excel file: pandas.read_excel()
 - Reading a MATLAB file: scipy.io.loadmat()
 - Reading shapefile and GeoJSON files: geopandas.read_file()
 - Reading GeoTIFF: rasterio.open()
 - Reading an image: matplotlib.pyplot.imread()



Data Cleaning

Data preparation: messy data → tidy data
Rectangular data structures → Data modelling
Intidy data:
• each variable forms a column

- each observation forms a row
- each cell is a single measurement





Wickham, H. (2014). Tidy Data. Journal of Statistical Software 59 (10). DOI: 10.18637/jss.v059.i10

From: https://www.openscapes.org/blog/2020/10/12/tidy-data/

UNIVERSITY OF TWENTE.

Exploratory Data Analysis (EDA)

EDA aims at **summarizing** the characteristics of a dataset with **statistical numbers** and **graphs**

Statistical Analysis + Visualization

Get an overview of the data

Orient further analysis \rightarrow choose correct methods/approaches

Help you to generate hypothesis

Spot problems in data

Understand properties of the variables (e.g., mean)

Understand relationships between variables



Statistics + Visualization

Anscombe's quartet							
I		П		Ш		IV	
x	у	x	у	x	у	x	у
10.0	8.04	10.0	9.14	10.0	7.46	8.0	6.58
8.0	6.95	8.0	8.14	8.0	6.77	8.0	5. 76
13.0	7.58	13.0	8.74	13.0	12.74	8.0	7.71
9.0	8.81	9.0	8.77	9.0	7.11	8.0	8.84
11.0	8.33	11.0	9.26	11.0	7.81	8.0	8.47
14.0	9.96	14.0	8.10	14.0	8.84	8.0	7.04
6.0	7.24	6.0	6.13	6.0	6.08	8.0	5.25
4.0	4.26	4.0	3.10	4.0	5.39	19.0	12.50
12.0	10.84	12.0	9.13	12.0	8.15	8.0	5.56
7.0	4.82	7.0	7.26	7.0	6.42	8.0	7.91
5.0	5.68	5.0	4.74	5.0	5.73	8.0	6.89

Property	Value
Mean of <i>x</i> in each case	9 (exact)
Variance of <i>x</i> in each case	11 (exact)
Mean of <i>y</i> in each case	7.50 (to 2 decimal places)
Variance of <i>y</i> in each case	4.122 or 4.127 (to 3 decimal places)
Correlation between x and y in each case	0.816 (to 3 decimal places)
Linear regression line in each case	y = 3.00 + 0.500x (to 2 and 3 decimal places, respectively)

UNIVERSITY OF TWENTE.

From: https://en.wikipedia.org/wiki/Anscombe%27s_quartet

Statistics + Visualization



EDA focuses on statistical analysis and visualization Visualization also can help to see the anomaly of statistical calculations Visualization can help to get specific insight if we work in spatiotemporal Tidy data is a standard way of mapping the meaning of a dataset to its structure Box plot can see the range of our data and we can see the outlier of the data Positive : separate the data No spatial correlation : random Negative : no correlation at all







Mean and Standard Deviation

Histogram and PDF

distribution of the data, showing the number of observations that fall within each bin. PDF is the continuous version of the histogram







Min, Max, Median, Percentile, Quartile

Percentile: Given a vector V of length N, the q-th percentile of V is the value q/100 of the way from the minimum to the maximum in a sorted copy of V.

Quartile: The q-th quantile of V is the value q of the way from the minimum to the maximum in a sorted copy of V.



Box plot: displays the five-number summary (the minimum, first quartile, median, third quartile, and maximum) of a set of data. It can tell you about your outliers and what their values are



https://towardsdatascience.com/understanding-boxplots-5e2df7bcbd51



Bar plots



From: https://matplotlib.org/





From: https://xkcd.com



Bi-Variate Analysis

Correlation

Relationship between two variables quantitatively

$$cor(x,y) = rac{cov(x,y)}{sd(x)sd(y)}$$
 $cov(x,y) = rac{1}{n}\sum_{i=1}^n (x_i-\overline{x})(y_i-\overline{y})$



Bi-Variate Analysis





Bi-Variate Analysis

Pair-plot





Exploratory Spatial Data Analysis

Geospatial data \rightarrow ESDA

"Traditional" EDA can be applied to spatial datasets for obtaining statistics and basic plots (barplot, histograms, boxplots,..).

ESDA tools connects a specific variable to a location/time It takes into account the values of the same variable in different locations/time.



Applying EDA to geospatial data



Observations Distributed Over Years



Spatial autocorrelation

Correlation of a variable with itself across space (in different places in space) \rightarrow relationships to neighbors

Positive spatial autocorrelation

values are similar to their neighbors or other close objects

clusters of similar values on the map

Zero or no spatial autocorrelation

random values of close objects or neighbors

no clear pattern visually

Negative spatial autocorrelation

values are dissimilar to their neighbors or close objects **dispersed** patterns of values on the map



Spatial autocorrelation

Positive spatial autocorrelation

No spatial autocorrelation

Negative spatial autocorrelation

From: (Radil, 2011)



Spatial autocorrelation



From: https://mgimond.github.io/Spatial/spatial-autocorrelation.html



SPATIAL AUTOCORRELATION: MORAN'S I

- *n* is the number of cases
- **x**_i is the variable value at a particular location
- **x**_j is the variable value at another location
- \overline{X} is the mean of the variable
- *w_{ij}* is a weight applied to the comparison between location *i* and location *j*

$$I = \frac{n \sum_{i} \sum_{j} w_{i,j} (x_{i} - \bar{x}) (x_{j} - \bar{x})}{\sum_{i} \sum_{j} w_{i,j} \sum_{i} (x_{i} - \bar{x})^{2}}$$



Check out the link below for more in-depth explanation: https://rpubs.com/corey_sparks/105700



Visualization on map





Connection map





From: https://www.data-to-viz.com/story/MapConnection.html

Box map







ESDA maps

Some examples of ESDA maps:

Box Map: <u>https://geodacenter.github.io/workbook/3a_mapping/lab3a.html#extreme-value-maps</u>

Brushing & linking:

https://www.spatialanalysisonline.com/HTML/eda esda and estda.htm

Conditional choropleth mapping:

http://publichealthintelligence.org/content/geography-diabetes-us-conditioned-map

Voronoi analysis: https://www.gislounge.com/voronoi-diagrams-and-gis/

Cartograms: https://gisgeography.com/cartogram-maps/

Connection map: <u>https://www.data-to-viz.com/story/MapConnection.html</u>



UNIVERSITY OF TWENTE.

Team Based Learning

Team based learning assignment

Ghelgheli decided to change his job, and as a tea lover, he opted to open a teahouse. He aimed to find the right location for his business, where many people were passing by and not many competitors around.

Ghelgheli started by collecting data, organizing it into rows and columns within a table on his computer. However, the data was somewhat messy, containing several missing values and even some anomalies. Nevertheless, Ghelgheli was enthusiastic about working with such a dataset. He used some cool techniques to clean the data, extract statistical measures, and generate plots and maps.

Through his analysis, Ghelgheli pinpointed a suitable location for his teahouse, and soon after opening, it became a local favorite.

Which data and methods do you think Ghelgheli utilized for his analysis? What interesting learnings did you derive from Ghelgheli's story? Can you provide some real-life examples similar to Ghelgheli's experience?



Data Collection: Ghelgheli started by collecting data on potential locations for his teahouse. This could include foot traffic data, competitor locations, rent prices, demographic information of the area, etc.

Data Cleaning: The data Ghelgheli collected was described as messy, with missing and strange values. Ghelgheli likely employed techniques like data imputation, outlier detection, and data validation to clean the dataset.

Statistical Analysis: Ghelgheli extracted statistical measures from the cleaned dataset. This could involve calculating means, medians, standard deviations, and other descriptive statistics to understand the characteristics of the data.

Visualization: Ghelgheli created plots and maps to visualize the data. This could include scatter plots, histograms, heatmaps, and geographical maps to identify patterns and trends in the data.



Decision Making: Through the analysis, Ghelgheli identified a suitable location for his teahouse based on the insights gained from the data analysis.

- The importance of data in decision-making processes
- The power of EDA techniques in uncovering insights and making informed decisions.
- How messy data can be transformed into valuable insights through proper cleaning and analysis.

