

Pyinterpolate: Spatial Interpolation in Python for point measurements and aggregated datasets

Szymon Moliński¹

DOI: [10.21105/joss.02869](https://doi.org/10.21105/joss.02869)

Software

- [Review](#) ↗
- [Repository](#) ↗
- [Archive](#) ↗

Editor: [Hugo Ledoux](#) ↗

Reviewers:

- [@chrisbrunsdon](#)
- [@kenohori](#)
- [@sdesabbata](#)

Submitted: 26 October 2020

Published: 28 June 2021

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

¹ Data Lions company, Poland, <https://datalions.eu>

Summary

Spatial Interpolation techniques are used to interpolate values at unknown locations and/or filter and smooth existing data sources. Those methods work for point observations and areal aggregates. The basic idea under this family of algorithms is that every point in a space can be described as a function of its neighbors values weighted by the relative distance from the analyzed point. It is known as the Tobler's First Law of Geography, which states: *everything is related to everything else, but near things are more related than distant things* (Tobler, 1970).

Kriging technique designed for mining applications exploits this statement formally and nowadays it has gained a lot of attention outside the initial area of interest. Today *kriging* is a set of methods which can be applied to problems from multiple fields: environmental science, hydrogeology, natural resources monitoring, remote sensing, epidemiology and ecology and even computer science (Chilès & Desassis, 2018). Most commonly Kriging interpolates values from point measurements or regular block units but many real-world datasets are different. Especially challenging are measurements of rates over areas of irregular shapes and sizes, as example administrative units in every country (Goovaerts, 2007).

Pyinterpolate is designed to tackle a problem of areas of irregular shapes and sizes with Area-to-Area and Area-to-Point Poisson Kriging functions. With those algorithms Pyinterpolate became an interpolation and filtering tool which is useful for social, environmental and public health sciences. Moreover, the package offers basic Kriging and Inverse Distance Weighting techniques and can be utilized in every field of research where geostatistical (distance) analysis gives meaningful results. Pyinterpolate merges basic Kriging techniques with more sophisticated Area-to-Area and Area-to-Point Poisson Kriging methods.

Statement of need

Pyinterpolate is a Python package for spatial interpolation and it is designed to perform predictions from point measurements and areal aggregates of different sizes and shapes. Pyinterpolate automates tasks performed by spatial statisticians, it helps with data exploration, semivariogram estimation and kriging predictions. Thing that makes Pyinterpolate different from other spatial interpolation packages is the ability to perform Kriging of areas of different shapes and sizes and this type of operation is extremely important in the context of social, medical and ecological sciences.

Importance of areal (block) Kriging

Areas of irregular shapes and sizes are especially challenging for analysis and modeling. The ability to transform areal aggregates into point support maps is desired by many applications.

39 As an example in public health studies data is aggregated over large areas due to the protection
40 of citizens' privacy but this process introduces bias to modeling and makes policy-making more
41 complex. The main three reasons behind transformation of choropleth maps with aggregated
42 counts into point support models are:

- 43 1. The presence of extreme unreliable rates that typically occur for sparsely populated areas
44 and rare events.
- 45 2. The visual bias resulting from aggregation of data over administrative units with various
46 shapes and sizes.
- 47 3. The mismatch of spatial supports for aggregated data and explanatory variables. This
48 prevents their direct use in models based on the correlation (Goovaerts, 2006).

49 In this context Area-to-Area Poisson Kriging serves as the noise-filtering algorithm or areal
50 interpolation model and Area-to-Point Poisson Kriging is designed to interpolate and transform
51 values and to preserve coherence of the prediction so the sum of average of disaggregated
52 estimates is equal to the baseline area value (Goovaerts & Gebreab, 2008). Area-to-Point
53 Poisson Kriging can be useful in the chained-models systems where change of support is
54 required to perform a study.

55 Researchers may use centroids of areas and perform point kriging over a prepared regular
56 point grid. However this method has pitfalls. Different sizes and shapes of units may lead
57 to imbalanced variogram point pairs per lag. Centroid-based approach does not catch spatial
58 variability of the linked variable, for example population density over area in the context of
59 infection rates.

60 To disaggregate areal data into point support one must know point support covariance and/or
61 semivariance of a regionalized variable. Then the semivariogram deconvolution is performed.
62 In this iterative process experimental semivariogram of areal data is transformed to fit the
63 semivariogram model of a linked point support variable. Example of it is the use of spatial
64 distribution of population to transform a semivariogram of disease rates which are the number
65 of cases divided by population. Semivariogram deconvolution is the core step of the Area-
66 to-Area and Area-to-Point Poisson Kriging operations. Poisson Kriging is widely used in the
67 social sciences, epidemiology and spatial statistics (Goovaerts, 2007; Goovaerts & Gebreab,
68 2008; Kerry et al., 2013).

69 Interpolation methods within Pyinterpolate

70 Package performs six types of spatial interpolation at the time of paper writing; five types of
71 Kriging and inverse distance weighting:

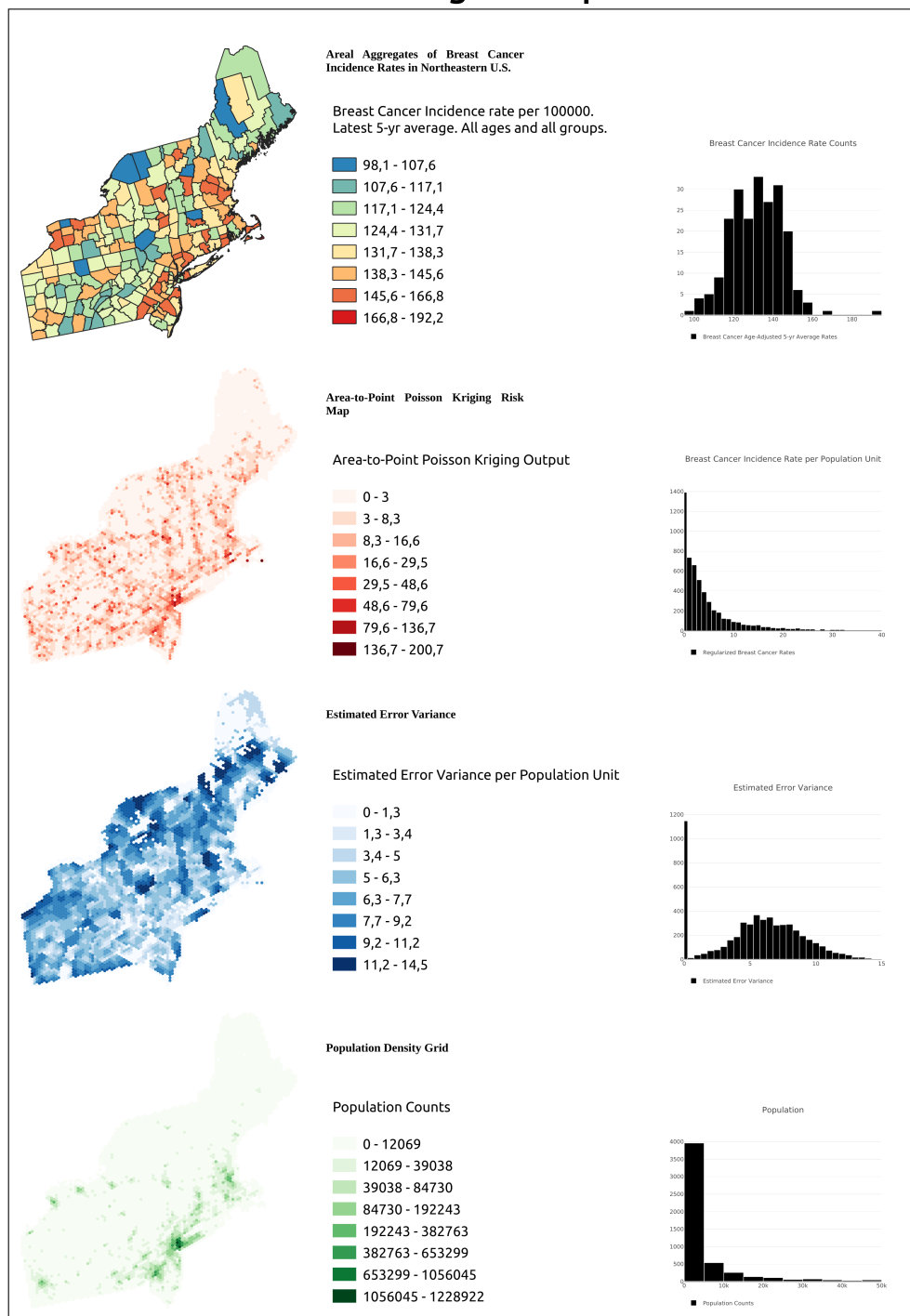
- 72 1. **Ordinary Kriging** which is a universal method for point interpolation.
- 73 2. **Simple Kriging** which is useful when the mean of the spatial process is known and it
74 is used for the point interpolation.
- 75 3. **Centroid-based Poisson Kriging**. This method of Kriging is based on the assumption
76 that each block can be collapsed into its centroid. It is much faster than Area-to-Area
77 and Area-to-Point Poisson Kriging but introduces bias related to the transformation of
78 areas into single points. It is used for areal interpolation and filtering.
- 79 4. **Area-to-Area Poisson Kriging**. Point support is included in the analysis and if it varies
80 over area. Model is able to catch this variation. It is used for areal interpolation and
81 filtering.
- 82 5. **Area-to-Point Poisson Kriging**. Areal support is deconvoluted in regards to the point
83 support. Output map has spatial resolution of the point support while coherence of anal-
84 ysis is preserved (sum of rates is equal to the output of Area-to-Area Poisson Kriging).
85 It is used for point-support interpolation and data filtering.

86 User starts with semivariogram exploration and modeling. Next researcher or algorithm
87 chooses the theoretical model which best fits the semivariogram. This model is used to
88 predict values at unknown locations. Areal data interpolation, especially transformation from
89 areal aggregates into point support maps, requires deconvolution of areal semivariogram. This
90 is an automatic process which can be performed without prior knowledge of kriging and spa-
91 tial statistics. The last step is Kriging itself. Poisson Kriging is especially useful for counts
92 over areas. On the other spectrum is Ordinary Kriging which is an universal technique which
93 works well with multiple point data sources. Predicted data is stored as a DataFrame known
94 from the *Pandas* and *GeoPandas* Python packages. Pyinterpolate allows users to transform
95 given point data into a regular numpy array grid for visualization purposes and to perform
96 large-scale comparison of different kriging techniques prediction output. Use case with the
97 whole scenario is available in the [paper package repository](#).

98 Package performs many steps automatically. User has the option to control prediction flow
99 with Python optional parameters in the function call. Package was initially developed for
100 epidemiological study, where areal aggregates of infections were transformed to point support
101 population-at-risk maps and multiple potential applications follow this algorithm. Initial field
102 of study (epidemiology) was the reason behind automation of many tasks related to data
103 modeling. It is assumed that users without a wide geostatistical background may use Pyinter-
104 polate for spatial data modeling and analysis, especially users which are observing processes
105 related to the human population.

106 The example of a process where deconvolution of areal counts and semivariogram regulariza-
107 tion occurs is presented in the [Figure 1](#).

Comparison of Real World data and Kriged Output



(C) Szymon Moliński, 2021

Figure 1: Structure of Pyinterpolate package.

108 Methodology

109 Chapter presents general methodology of calculations within package. Concrete use case is
110 presented in document [here](#). Comparison of algorithm with the **gstat** package is available
111 [here](#).

112 Spatial Interpolation with Kriging

113 Kriging, which is the baseline of the Pyinterpolate package, is an estimation method that gives
114 the best unbiased linear estimates of point values or block averages ([Armstrong, 1998](#)). Kriging
115 minimizes variance of a dataset with missing values. Baseline technique is the **Ordinary**
116 **Kriging** where value at unknown location \hat{z} is estimated as a linear combination of K neighbors
117 with value z and weights λ assigned to those neighbors (1).

(1)

$$\hat{z} = \sum_{i=1}^K \lambda_i z_i$$

118 Weights λ are a solution of following system of linear equations (2):

(2)

$$\sum_{j=1}^K \lambda_j C(x_i, x_j) - \mu = \bar{C}(x_i, V); i = 1, 2, \dots, K$$

$$\sum_i \lambda_i = 1$$

119
120 where $C(x_i, x_j)$ is a covariance between points x_i and x_j , $\bar{C}(x_i, V)$ is an average covariance
121 between point x_i and all other points in a group (K points) and μ is a process mean. The
122 same system may be solved with semivariance instead of covariance (3):

(3)

$$\sum_{j=1}^K \lambda_j \gamma(x_i, x_j) + \mu = \bar{\gamma}(x_i, V); i = 1, 2, \dots, K$$

$$\sum_i \lambda_i = 1$$

123
124 where $\gamma(x_i, x_j)$ is a semivariance between points x_i and x_j , $\bar{\gamma}(x_i, V)$ is an average semi-
125 variance between point x_i and all other points. Semivariance is a key concept of spatial
126 interpolation. It is a measure of a dissimilarity between observations in a function of distance.
127 Equation (4) is a experimental semivariogram estimation formula.

(4)

$$\frac{1}{2N} \sum_i^N (z_{(x_i+h)} - z_{x_i})^2$$

where z_{x_i} is a value at location x_i and $z_{(x_i+h)}$ is a value at translated location in a distance h from x_i .

In the next step theoretical models are fitted to the experimental curve. Pyinterpolate package implements linear, spherical, exponential and gaussian models but many others are applied for specific cases (Armstrong, 1998). Model with the lowest error is used in (3) to estimate γ parameter.

Ordinary Kriging is one of the classic Kriging types within the package. **Simple Kriging** is another available method for point interpolation. Simple Kriging may be used when the process mean is known over the whole sampling area. This situation rarely occurs in real world. It can be observed in places where sampling density is high (Armstrong, 1998). Simple Kriging system is defined as:

$$(5) \quad \hat{z} = R + \mu$$

where μ is a process mean and R is a residual at a specific location. Residual value is derived as the first element (denoted as 1) from:

$$(6) \quad R = ((Z - \mu) \times \lambda) \mathbf{1}$$

Number of values depends on the number of neighbours in a search radius, similar to equation (1) for Ordinary Kriging. λ weights are the solution of following function:

$$(7) \quad \lambda = K^{-1}(\hat{k})$$

where K is a semivariance matrix between each neighbour of size $N \times N$ and k is a semivariance between unknown point and known points of size $N \times 1$.

Package allows use of the three main types of Poisson Kriging: Centroid-based Poisson Kriging, Area-to-Area Poisson Kriging and Area-to-Point Poisson Kriging. Risk over areas (or points) for each type of Poisson Kriging is defined similarly to the equation (1) but weights associated with the λ parameter are estimated with additional constraints related to the population weighting. The spatial support of each unit needs to be accounted for in both the semivariogram inference and kriging. Full process of areal data Poisson Kriging is presented in (Goovaerts, 2006) and semivariogram deconvolution which is an intermediate step in Poisson Kriging is described in (Goovaerts, 2007).

Modules

Pyinterpolate is designed from seven modules and they cover all operations needed to perform spatial interpolation: from input/output operations, data processing and transformation, semivariogram fit to kriging interpolation. Figure 2 shows package structure.

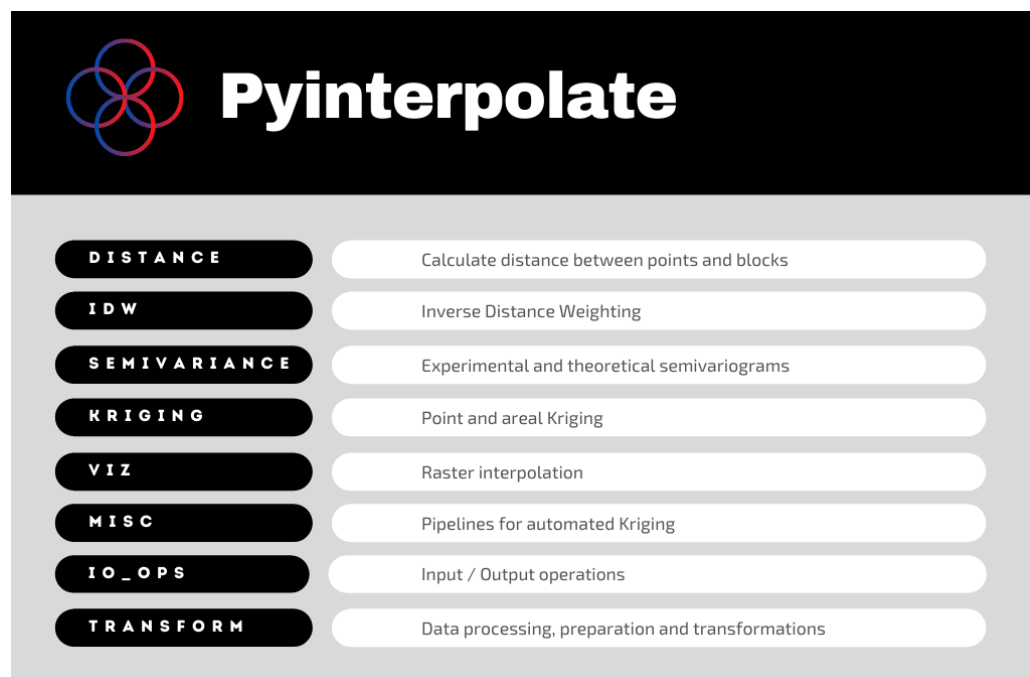


Figure 2: Structure of Pyinterpolate package.

Modules follow typical data processing and modeling steps. The first module is **io_ops** which reads point data from text files and areal or point data from shapefiles, then changes data structure for further processing. **Transform** module is responsible for all tasks related to changes in data structure during program execution. Sample tasks are:

- finding centroids of areal data,
- building masks of points within lag.

Functions for distance calculation between points and between areas (blocks) are grouped within **distance** module. **Semivariance** is most complex part of Pyinterpolate package. It has three special classes for calculation and storage of different types of semivariograms (experimental, theoretical, areal and point types). **Semivariance** module has other functions important for spatial analysis:

- function for experimental semivariance / covariance calculation,
- weighted semivariance estimation,
- variogram cloud preparation,
- outliers removal.

Kriging module contains three main types of models Ordinary and Simple Kriging models as well Poisson Kriging of areal counts models. Areal models are derived from (Goovaerts & Gebreab, 2008), simple Kriging and ordinary Kriging models are based on (Armstrong, 1998).

It is possible to show output as numpy array with **viz** module and to compare multiple kriging models on the same dataset with **misc** module. Evaluation metric for comparison is an average root mean squared error over multiple random divisions of a passed dataset.

178 Comparison to Existing Software

179 Pyinterpolate is one package from a large ecosystem of spatial modeling and spatial inter-
180 polation packages written in Python. The main difference between Pyinterpolate and other
181 packages is focus on areal deconvolution methods and Poisson Kriging techniques useful for
182 ecology, social science and public health studies in the presented package. Potential users
183 may choose other packages if their study is limited to point data interpolation.

184 The most similar and most important package from Python environment is **PyKrige** (Murphy
185 et al., 2020). PyKrige is designed especially for point kriging. PyKrige supports 2D and 3D
186 ordinary and universal Kriging. User is able to incorporate own semivariogram models and/or
187 use external functions (as example from **scikit-learn** package (Pedregosa et al., 2011)) to
188 model drift in universal Kriging. Package is well designed, and it is actively maintained.

189 **GRASS GIS** (GRASS Development Team, 2020) is well-established software for vector and
190 raster data processing and analysis. GRASS contains multiple modules and GRASS function-
191 alities can be accessed from multiple interfaces: GUI, command line, C API, Python API,
192 Jupyter Notebooks, web, QGIS and R. GRASS has two functions for spatial interpolation:
193 `r.surf.idw` and `v.surf.idw`. Both use Inverse Distance Weighting technique, first interpo-
194 lated raster files and second vectors (points).

195 **PySAL** is next GIS / geospatial package which can be used to interpolate missing values – but
196 this time at areal scale. Package's **tobler** module can be used to interpolate areal values of
197 specific variable at different scales and sizes of support (knaap et al., 2020). Moreover, package
198 has functions for multisource regression, where raster data is used as auxiliary information to
199 enhance interpolation results. Conceptually tobler package is close to the Pyinterpolate, where
200 main algorithm transforms areal data into point support derived from auxiliary variable.

201 **R programming language** offers **gstat** package for spatial interpolation and spatial mod-
202 eling (Pebesma, 2004). Package is designed for variogram modelling, simple, ordinary and
203 universal point or block kriging (with drift), spatio-temporal kriging and sequential Gaussian
204 (co)simulation. Gstat is a solid package for Kriging and spatial interpolation and has the
205 largest number of methods to perform spatial modelling. The main difference between gstat
206 and Pyinterpolate is availability of area-to-point Poisson Kriging based on the algorithm pro-
207 posed by Goovaerts (Goovaerts, 2007) in Pyinterpolate package. Comparison to **gstat** is
208 available in the [paper repository](#).

209 Appendix

- 210 1. [Paper repository with additional materials](#)
- 211 2. [Package repository](#)

212 References

- 213 Armstrong, M. (1998). *Basic linear geostatistics*. Springer. [https://doi.org/10.1007/](https://doi.org/10.1007/978-3-642-58727-6)
214 [978-3-642-58727-6](https://doi.org/10.1007/978-3-642-58727-6)
- 215 Chilès, J.-P., & Desassis, N. (2018). Fifty years of kriging. In B. S. Daya Sagar, Q. Cheng, & F.
216 Agterberg (Eds.), *Handbook of mathematical geosciences: Fifty years of IAMG* (pp. 589–
217 612). Springer International Publishing. https://doi.org/10.1007/978-3-319-78999-6_29
- 218 Goovaerts, P. (2006). Geostatistical analysis of disease data: Accounting for spatial support
219 and population density in the isopleth mapping of cancer mortality risk using area-to-point

- 220 poisson kriging. *International Journal of Health Geographics*, 5. <https://doi.org/10.1186/1476-072X-5-52>
- 221
- 222 Goovaerts, P. (2007). Kriging and semivariogram deconvolution in the presence of irregular
- 223 geographical units. *Mathematical Geosciences*, 40, 101–128. <https://doi.org/10.1007/s11004-007-9129-1>
- 224
- 225 Goovaerts, P., & Gebreab, S. (2008). How does poisson kriging compare to the popular
- 226 BYM model for mapping disease risks? *International Journal of Health Geographics*, 7, 6. <https://doi.org/10.1186/1476-072x-7-6>
- 227
- 228 GRASS Development Team. (2020). *Geographic resources analysis support system (GRASS GIS) software*. Open Source Geospatial Foundation. <https://grass.osgeo.org>
- 229
- 230 Kerry, R., Goovaerts, P., Smit, I. P. J., & Ingram, B. R. (2013). A comparison of multiple in-
- 231 dicator kriging and area-to-point poisson kriging for mapping patterns of herbivore species
- 232 abundance in kruger national park, south africa. *International Journal of Geographical*
- 233 *Information Science*, 27, 47–67. <https://doi.org/10.1080/13658816.2012.663917>
- 234 knaap, eli, Cortes, R. X., Rey, S., Gaboardi, J., & Frontiera, P. (2020). *Pysal/tobler: Release*
- 235 *v0.5.4* (Version v0.5.4) [Computer software]. Zenodo. <https://doi.org/10.5281/zenodo.4385980>
- 236
- 237 Murphy, B., Müller, S., & Yurchak, R. (2020). *GeoStat-framework/PyKrige v1.5.1* (Version
- 238 v1.5.1) [Computer software]. Zenodo. <https://doi.org/10.5281/zenodo.3991907>
- 239
- 240 Pebesma, E. J. (2004). Multivariable geostatistics in s: The gstat package. *Computers &*
- 241 *Geosciences*, 30(7), 683–691. <https://doi.org/10.1016/j.cageo.2004.03.012>
- 242
- 242 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel,
- 243 M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau,
- 244 D., Brucher, M., Perrot, M., & Duchesnay, E. (2011). Scikit-learn: Machine learning in
- 245 Python. *Journal of Machine Learning Research*, 12, 2825–2830.
- 246
- 246 Tobler, W. R. (1970). A computer movie simulating urban growth in the detroit region.
- 247 *Economic Geography*, 46, 234. <https://doi.org/10.2307/143141>