

Spatial Signatures for Geographic Feature Types: Examining Gazetteer Ontologies using Spatial Statistics

Rui Zhu, Yingjie Hu, Krzysztof Janowicz and Grant McKenzie STKO Lab

Department of Geography, University of California, Santa Barbara





Outline



- Motivations
- Datasets
- Methods
- Case Studies
- Future Work



Motivations



• A wide variety of gazetteers



July 28th 2016

Esri User Conference 2016



- However, these gazetteers differ largely in terms of:
 - Overall coverage
 - Underlying data sources
 - Provided functionalities
 - Geographic feature type ontologies





• Examples:

Dam





Stream

GeoNames

TGN











July 28th 2016



• Examples



July 28th 2016





 How to understand such semantic heterogeneity among gazetteers?



Janowicz, K., 2012. Observation-driven geo-ontology engineering. Transactions in GIS 16 (3), 351–374.

July 28th 2016



Datasets





DBpedia Places





GeoNames

Getty Thesaurus of Geographic Names (TGN)

Extracted from DBpedia articles Formal geographical database that contains over eight million place names

Focus on places that are culturally or historically significant

73 feature types in US

198 feature types in US

285 feature types in US

July 28th 2016



Datasets (cont.)



- Common information
 - Toponyms/Place names
 - Geographic feature type
 - Spatial footprints



Dams in GeoNames



Methods





July 28th 2016





July 28th 2016

Esri User Conference 2016



• Spatial point pattern analysis (cont.)

spatial statistics (local)

- Intensity of point patterns
- Distance to nearest neighbor
- Ripley's K (i.e. range and mean deviation from the theoretical values)
- Kernel density estimation (i.e. bandwidth and range)
- Standard deviational ellipse (i.e. rotation, std. along x-axis and y-axis)





1.4e-08

09 4e-09 6e-09 8e-09 1e-08

• Spatial point pattern analysis (cont.)

spatial statistics (global)

- Overall intensity of point patterns
- Kernel density estimation (i.e. bandwidth and range)

Dams in GeoNames









• Spatial autocorrelation analysis Conversion: point data → raster map

Dams in GeoNames



Cell size: 36km * 22.2 km

Cell value: number of places falling in the cell

July 28th 2016





• Spatial autocorrelation analysis (cont.)

Spatial statistics

- Global Moran's I
- Sample Semivariogram (i.e. semivariances at first, median and last lag distances).



Experimental Semivariogram for Geonames Dam

July 28th 2016





Spatial interaction with other geographic features

Spatial statistics (Internal)

Count of distinct nearest feature types

 $CountNearest_i = \frac{m_i}{N}$

where m_i is the number of distinct nearest feature type for feature type *i* and *N* is the total number of feature types in one gazetteer.

• Entropy of nearest feature types

EntropyNearest_i =
$$- \overset{N}{\bigcirc} \frac{n_j}{j=1} \frac{n_j}{N} \log(\frac{n_j}{N})$$

where N is the total number of feature types in one gazetteer and n_j is the number of nearest instances from the j^{th} feature type.

July 28th 2016





• Spatial interaction with other geographic features

Spatial statistics (External)

Population (LandScan2014) Population for each feature point Minimum • Maximum • Mean Standard deviation

July 28th 2016







July 28th 2016

Esri User Conference 2016



Overview

- Same name and similar spatial patterns
- Same name but different spatial patterns
- Different names but similar spatial patterns
- Different names and different spatial patterns





• Same name and similar spatial patterns



July 28th 2016

Coordinate 1



• Same name but different spatial patterns



MDS (2D) for Mountain

July 28th 2016

Coordinate 1



• Different names but similar spatial patterns



MDS (2D) for AdministrativeRegion (DBpedia Places), ADM2(GeoNames), County (TGN)

July 28th 2016

Coordinate 1



• Different names and different spatial patterns



MDS (2D) for DBpeida Places

July 28th 2016

Future work



- Derive additional statistical features to represent the spatial signature
 - (e.g. statistics for co-occurrence, topological relations);
- Quantify the dissimilarity/similarity of place types using such spatial signatures
 - (e.g. supervised/ unsupervised learning algorithms);
- Combine spatial signatures with previously studied temporal and thematic signatures;
- Integrate this study (bottom-up) with classical top-down knowledge engineering.





Thank you! Questions and comments?









- *. MeanNearestDistance: Mean distance to nearest neighbor
- *. VarNearestDistance: Variance distance to nearest neighbor
- *. LocalIntensity: Local intensity
- *. RipleyKRange: Ripley's K (range)
- *. RipleyKMeanDev: Ripley's K (mean deviation)
- *. LocalKernelBW: Local kernel density (bandwidth)
- *. LocalKernelRange: Local kernel density (range)

- *. EllipseYStdDev: Standard deviational ellipse (std dev along x-axis) *. PopMax: Population value (max)
- *. Moranl: Global Moran's I
- *. FirstSemivar: Semivarigram value (at first distance lag)
- *. MedianSemivar: Semivariogram value (at median distance lag)
- *. LastSemivar: Semivariogram value (at last distance lag)
- *. GlobalIntensity: Global Intensity
- *. GlobalKernelBW: Global kernel density (bandwidth)

- *. PopMean: Population value (mean)
- *. PopSD: Population value (std dev)
- *. RoadMin: Shortest distance to road (min)
- *. RoadMax: Shortest distance to road (max)
- *. RoadMean: Shortest distance to road (mean)
- *. RoadSD: Shortest distance to road (std dev)

July 28th 2016





| Spatial Signatures | |
|----------------------------------------------------|-----------------------------------|
| Mean distance to nearest neighbor | Global kernel density (bandwidth) |
| Variance distance to nearest neighbor | Population value (min) |
| Local intensity | Population value (max) |
| Ripley's K (range) | Population value (mean) |
| Ripley's K (mean deviation) | Population value (std dev) |
| Local kernel density (bandwidth) | Short distance to road (min) |
| Standard deviational ellipse (rotation) | Short distance to road (max) |
| Standard deviational ellipse (std dev along y-axis | Short distance to road (std dev) |
| Semivariogram (median distance lag) | Entropy of nearest feature type |

