

Understanding the Semantics of Places in Gazetteers via Spatial Analysis

Rui Zhu

STKO Lab Department of Geography University of California, Santa Barbara



Rui Zhu

Outline



- Motivations
- Data sets
- Methods
- Results & Discussions
- Future work





Contraction of the second seco

Motivations

• A wide variety of gazetteers





Rui Zhu

Motivations (cont.)



• Traditional techniques for integration/ alignment:

- Expert guess
- String similarity measures; e.g. Levenshtein distance
- Network similarity measures. e.g. Structure equivalence





Motivations (cont.)



• However, their geo-ontologies/ typing schema are different!







Geography USB Geography Solution

Motivations (cont.)

How to understand such semantic heterogeneity among gazetteers?



AAG 2016

Rui Zhu

Data Sets









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

Number of feature types

72

234

285



Rui Zhu

Geography CSB Geography

Data Sets (cont.)

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



Dams in GeoNames







Methods



• Overview:









• Spatial Point Pattern Analysis Sampling for local analysis









- Spatial Point Pattern Analysis (cont.)
 - Spatial Semantic Signatures (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 deviation ellipse (i.e. rotation, std. along x-axis and y-axis)







Ripley's K, Kernel density estimateion and Standard deviation ellipse for Dams in GeoNames



- Spatial Point Pattern Analysis (cont.) Spatial Semantic Signatures (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 : 36 km * 22.2 km

Cell value: number of instances falling in the cell



- Spatial Autocorrelation Analysis (cont.) Spatial Semantic Signatures
 - Global Moran's I
 - Sample Semivariogram (i.e. semivariances at first, median and last lag distances).



AAG 2016



• Spatial Interaction with other geographic features

Spatial Semantic Signatures

Population (LandScan2014)

Population for each feature point

Population Value Road Segment (Digital Chart of the World)

Distance to nearest segment for each feature point









• Same names and similar spatial patterns



MDS (2D) for Park



Results and Discussions (cont.)

• Same names but different spatial pattern



MDS(2D) for Mountain

AAG 2016

Spatial Data Mining & Big Data Analytics

UCSB

geography

Results and Discussions (cont.)

• Different names but similar spatial pattern

MDS (2D) for Administrative (DBpedia Places), ADM2(GeoNames), County (TGN) 9 • DBpedia Places GeoNames • TGN S Coordinate 2 0 ĥ -10 -5 0 Coordinate 1

AAG 2016

UCSB

geography

Results and Discussions (cont.)

• Different names and different spatial pattern



MDS (2D) for different types in DBpedia Places



Spatial Data Mining & Big Data Analytics

UCSB

geography

Future work



- Derive additional statistical features to represent the spatial semantic signature (e.g. statistics for co-occurrence, topological relations);
- Quantify the dissimilarity/similarity of place types using such spatial semantic 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.





Special thanks to:

Yingjie Hu, Krzysztof Janowicz and Grant McKenzie

Questions and/or comments?



