



# Understanding Semantic Uncertainty in Volunteered Geographic Information

Rui Zhu

University of California, Santa Barbara

AAG 2019 Session: GIScience in the Post-Truth Era I



# Outline

- Motivations
- Methodology
- Experiments
- Summary & Discussions



# Motivations





# Motivation

Uncertainty of VGI:

- Positions: coordinates, addresses, postcodes, ...
- Geometric representations and topological relations
- Attributes: number of check-ins, number of employees, ...
- **Semantics: place types**



# Motivation

## Semantic Uncertainty of VGI

- *Restaurants* in Foursquare and Google places might not be the same.
- *Mountains* in DBpedia Places are different from *Mountains* in GeoNames.
- A place should be labeled more likely as a restaurant or as a bar?
- Will spatial contexts help to reduce such uncertainties?

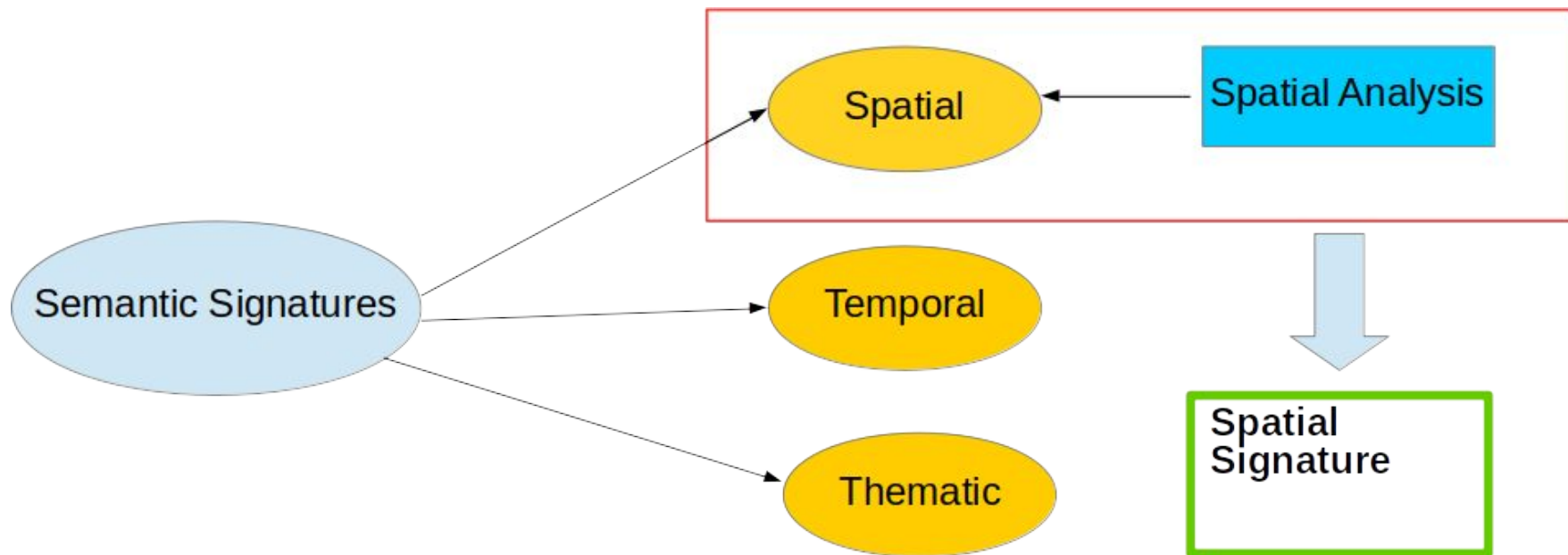


# Methodology

## Conventional approach for uncertainty analysis

- Numeric data
  - Distribution of the data → variance, entropy, Bayes theorem
- Categorical data
  - Indicator statistics
  - **Semantic signatures** → semantic distribution of place types

# Methodology



Janowicz, K., McKenzie, G., Hu, Y., Zhu, R., and Gao, So. (2018): [Using Semantic Signatures for Social Sensing in Urban Environments](#). Mobility Patterns, Big Data and Transport Analytics.



# Methodology

## Spatial Signatures

- ***Spatial structure*** of the data belonging to a place type is used to quantify its semantics.
- ***Spatial statistics*** are applied to describe such spatial structure.
- Spatial point patterns, Spatial autocorrelation analysis, spatial interaction analysis with other geographic features, place-based analysis. → **41 statistics**



# Methodology

## Spatial Signatures - Spatial point patterns

- Intensity-based: local intensity, kernel density estimation
- Distance-based: nearest-neighbor distance, Ripley's K, and standard deviational analysis

Randomly Selected Points using CSR in Contiguous US



Generate random points  
(Complete Spatial Randomness)

Geonames Dens in Contiguous US



Spatial Point Pattern for the 488 Randomly Selected Sample  
(Geonames Dens)



Select nearest 100 neighbors  
for each random points



Kernel Density Map of Spatial Point Pattern for the 488 Randomly Selected Sample



Average the statistics  
over all random points

Standard Deviation Ellipse of Spatial Point Pattern  
for the 488 Randomly Selected Sample (Geonames Dens)



Conduct spatial point pattern  
analysis on these 100 neighbors

# Methodology

## Spatial Signatures - Spatial point patterns - Examples

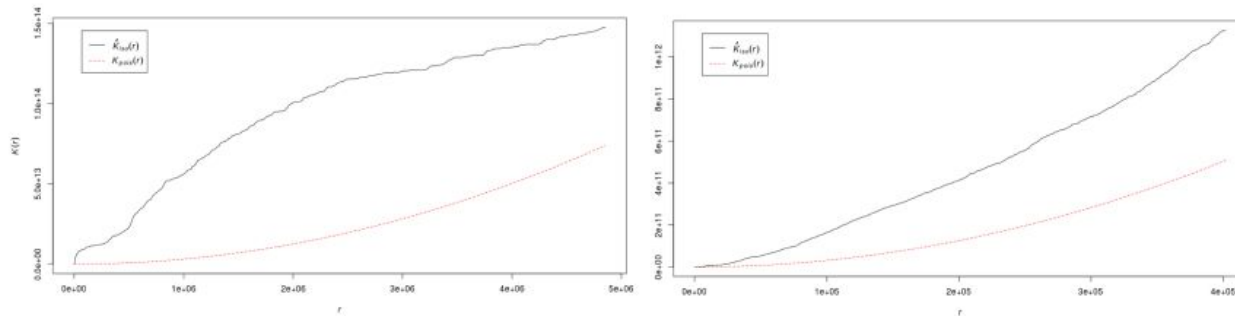


Figure 1: Ripley's K of *Park* (left) and *Dam* (right) from DBpedia Places.

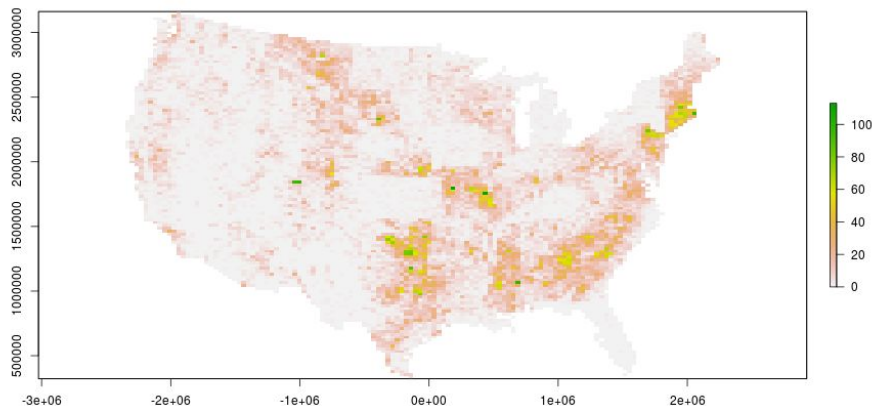
Statistics: mean and std. of the deviation between theoretical can observed K curves

# Methodology

## Spatial Signatures - Spatial Autocorrelation Analysis

- Moran's I: how intensities of cells differ from their neighbors
- Semivariogram: measure the variation of cell intensities in a specific distance lag class.

Dams in GeoNames



Cell size : 36 km \* 22.2 km

Cell value: number of  
instances falling in the cell

# Methodology

## Spatial Signatures - Spatial Autocorrelation Analysis - Examples

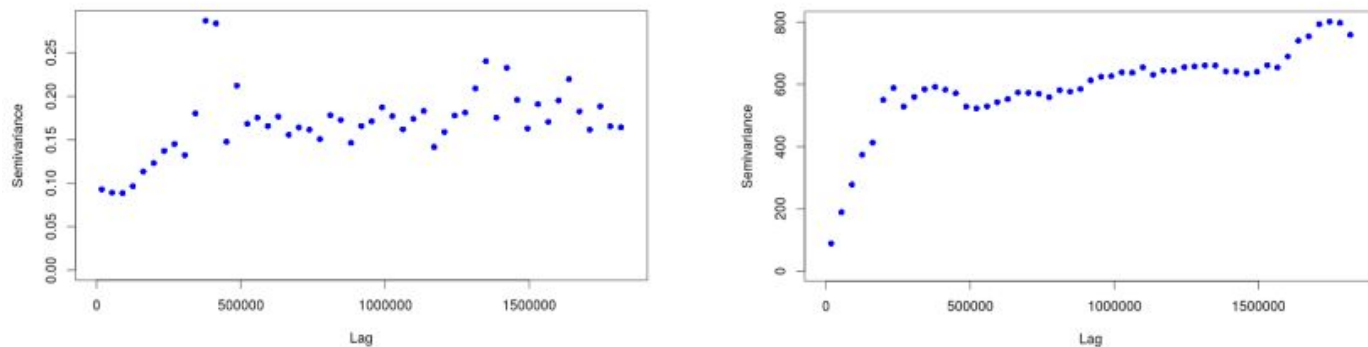


Figure 2: Experimental semivariogram of *Park* (left) and *Dam* (right) from TGN.

Statistics: mean and std. of the semivariance at first, median and last lag distance

# Methodology

## Spatial Signatures - Spatial Interaction with Other Geographic features

- Population
- Climate
- Road network

Population (LandScan2014)

Population for each feature point



Road Segment (Digital Chart of the World)

Distance to nearest segment for each feature point

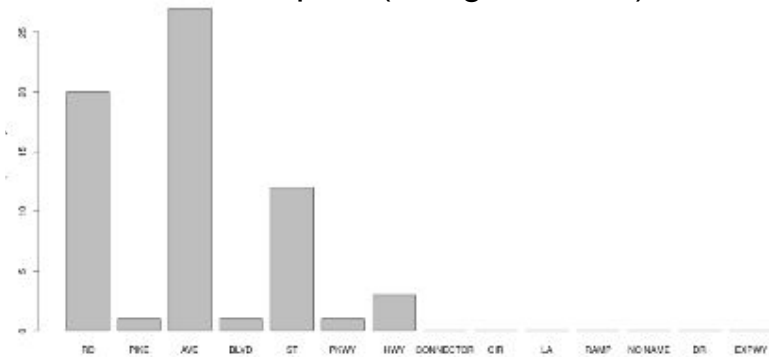


- Minimum
- Maximum
- Mean
- Standard deviation

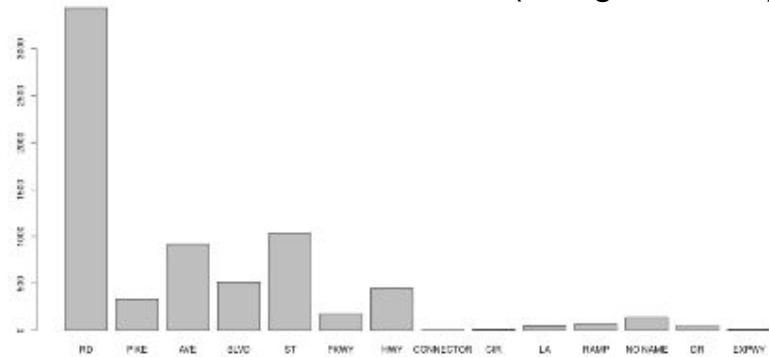
# Methodology

## Spatial Signatures - Spatial Interaction with Other Geographic features - Examples

amusement park (Google Places)



restaurant (Google Places)

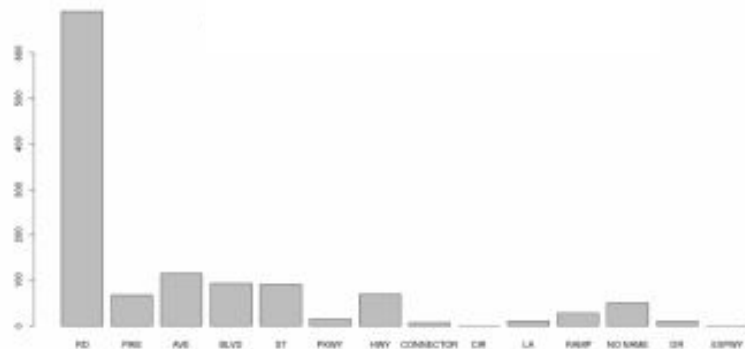


**Road suffix distribution of *amusement park* and *restaurant* from Google Places**

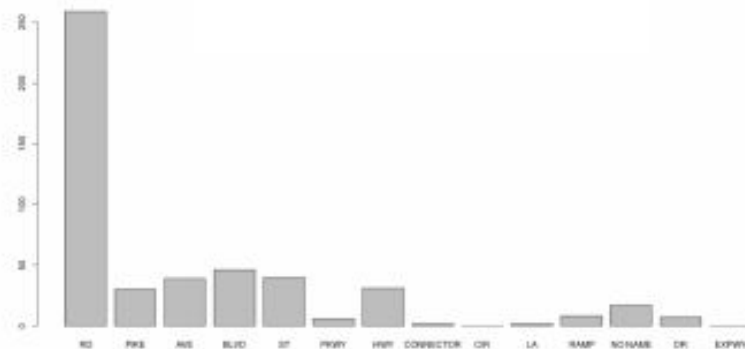
# Methodology

## Spatial Signatures - Spatial Interaction with Other Geographic features - Examples

car dealer (Google Places)



car dealership (Foursquare)



**Road suffix distribution** of *car dealer* from Google Places and *car dealership* from Foursquare



# Methodology

## Spatial Signatures - Place-based statistics

In contrast to spatial statistics, place-based statistics focus more on describing the ***topological*** and ***hierarchical relations*** between places.

- The number (and entropy) of distinct states (or counties) a place type occurs in;
- The number (and entropy) of adjacent states (or counties) that also contain places of the same type;





# Methodology

## Spatial Signatures - Place-based statistics - Examples

- To distinguish feature types:
  - Glacier: found in eight US-states according to DBpedia
  - River: found in all states



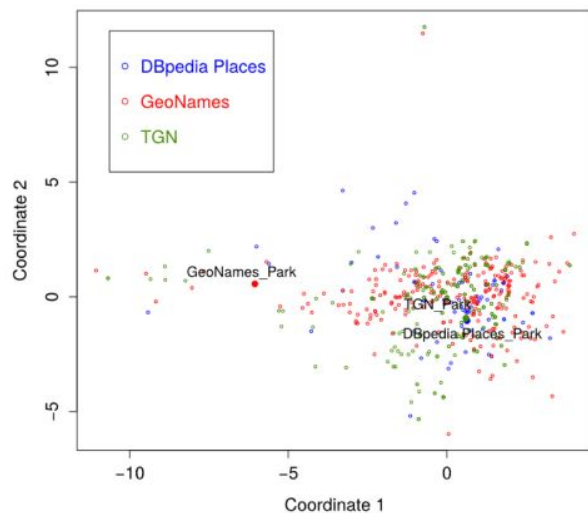
# Methodology

Spatial Point Pattern		Spatial Autocorrelations	Spatial Interaction with Other Geographic Features		Place-based statistics	
Local	Intensity	Global Moran's $I$	Population	min	Number of distinct states (or counties)	
	Mean distance to nearest neighbor			max		
	std. of distance to nearest neighbor			mean	Entropy of states (or counties)	
	Kernel density (range)			std.	Number of adjacent states (or counties) that have the same feature type	
	Kernel density (bandwidth)		Road Network	min of shortest distance		
	Ripley's $K$ (range)			max of shortest distance		
	Ripley's $K$ (mean deviation)			mean of shortest distance	Number of distinct feature types for nearest neighbor	
	std. ellipse (rotation)	std. of shortest distance				
	std. ellipse (std. along x-axis)	Semivariogram (first distance lag)	entropy of nearest road types	Entropy of feature types for nearest neighbor		
	std. ellipse (std. along y-axis)		mean precipitation			
Global	Intensity		Semivariogram (median distance lag)	Climate	std. precipitation	LDA-based approach
		mean temperature max				
	std. temperature max					
	mean temperature min					
	Kernel density (range)	Semivariogram (last distance lag)	std. temperature min			Entropy of the topic distribution
			mean water vapor pressure			
	Kernel density (bandwidth)		std. water vapor pressure			

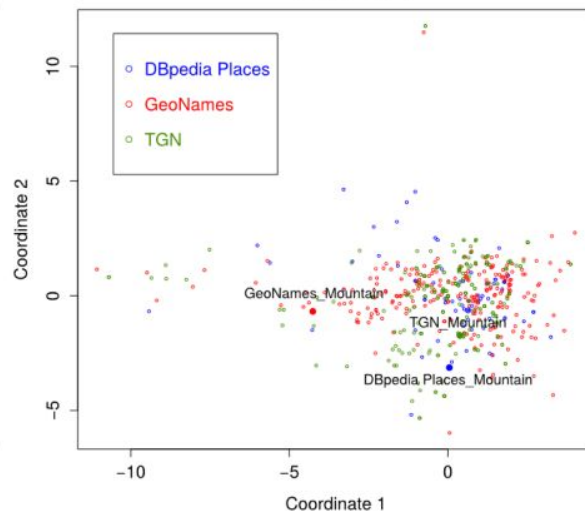
# Experiments

## 1. Similarity of place types

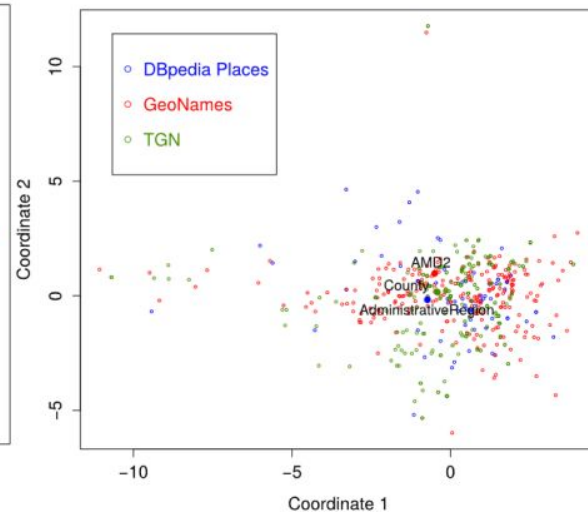
Case 1



Case 2



Case 3




## 2. Coreference resolution



# Which Kobani?

Have to use the feature types in addition to string and spatial distances.


Dedipedia

[Browse using ↻](#)
[Formats ↻](#)
[Faceted Browser ↻](#)
[Spans! ↻](#)

## About: Kobani

An Entry of Type: [settlement](#), from Named Graph: [http://dbpedia.org](#), within Data Space: [dbpedia.org](#)

Kobani (Kurdish: کوبانی pronounced [koˈbaːni], also rendered Kobanê [koˈbaːne], also known as Ayn al-Arab (Arabic: عين العرب North Levantine pronunciation: [ʕeːn elˈʕarab]), is a city in the Aleppo Governorate in northern Syria, lying immediately south of the border with Turkey. As a consequence of the Syrian Civil War, the city has been under control of the Kurdish YPG militia since 2012.

Property	Value
<span>dbpedia:PopulatedPlace/areaTotal</span>	<span>7.0</span>
<span>dbpedia:abstract</span>	<p><span>Kobani</span> (Kurdish: <span>کوبانی</span> pronounced [koˈbaːni], also rendered Kobanê [koˈbaːne], also known as Ayn al-Arab (Arabic: <span>عين العرب</span> North Levantine pronunciation: [ʕeːn elˈʕarab]), is a city in the Aleppo Governorate in northern Syria, lying immediately south of the border with Turkey. As a consequence of the Syrian Civil War, the city has been under control of the Kurdish YPG militia since 2012. In 2014, it was unofficially declared to be the administrative center of the Kobani Canton of Rojava. From September 2014 to January 2015, the city was under siege by Islamic State of Iraq and the Levant. Most of the city was destroyed and most of the population fled to Turkey. In 2015, many returned and reconstruction began. Prior to the Syrian Civil War, Kobani was recorded as having a population of close to 45,000. The majority of inhabitants were Kurds, with Arab, Turkmen, and Armenian minorities. <span>ⓘ</span></p>
<span>dbpedia:areaTotal</span>	<span>7000000.000000</span> (xsd:double)
<span>dbpedia:country</span>	<span>dbpedia:Syria</span>
<span>dbpedia:elevation</span>	<span>520.000000</span> (xsd:double)
<span>dbpedia:partOf</span>	<span>dbpedia:Ayn_al-Arab_District</span> <span>dbpedia:Aleppo_Governorate</span>



# Summary & Discussions

- Semantic uncertainty of VGI has to be understood and quantified
- Semantic signatures are introduced to quantify the semantic uncertainty

In the future:

- Need a framework/guideline of using semantic signatures
- To use semantic uncertainty to infer other types of uncertainties
- From exploratory study to solving emergent VGI challenges:
  - federated geographic information retrieval, place alignment, data cleaning, ...
- More advanced spatial /patial statistics could be incorporated into the signature set



**Thanks a lot!**

**Any questions / comments?**