Discovering Spatial and Temporal Links among RDF Graphs

Publishing and Interlinking Linked Geospatial Data In Conjunction with the 12th Extended Semantic Web Conference

Portoroz, Slovenia, 1st June 2015



LE Meloies Univ

HELLENIC REPUBLIC National and Kapodistrian University of Athens

Presenter: Panayiotis Smeros





- Introduction to Entity Resolution and Link Discovery
 - Examples, Definitions, Common Problems
- Spatial Entity Resolution
- Spatial and Temporal Link Discovery
 - Background and Developed Methods
 - Extensions to the Silk Framework
 - Hands-on



Entities in Real-World



Most of our knowledge about the world is based on **entities** and their **relations**:







Many names, descriptions or IDs (URIs) are used for the same **real-world entity**:







Many applications provide valuable information about each of these **entities**:







Many applications provide valuable information about each of these **entities**:







Problem of understanding that two (or more) entities in **data-world** are references of the same **real-world** entity. [Christen, TKDE'11]





Problem of understanding that two (or more) entities in **data-world** are references of the same **real-world** entity. [Christen, TKDE'11]



Discovering Spatial and Temporal Links among RDF Graphs





Problem of understanding that two (or more) entities in **data-world** are references of the same **real-world** entity. [Christen, TKDE'11]



Discovering Spatial and Temporal Links among RDF Graphs

Discovering Spatial and Temporal Links among RDF Graphs

10

Problem of understanding that two (or more) entities in **data-world** are references of the same **real-world** entity. [Christen, TKDE'11]











Let *S* and *T* be two sets of entities. We define a distance (similarity) function $d_{similarity}$ and a distance (similarity) threshold $\theta_{d_{similarity}}$ as follows:

$$d_{similarity} \colon S \times \mathbf{T} \to [0,1]$$
 , $\theta_{d_{similarity}} \in [0,1]$

We define the set of discovered similarity links $DL_{similarity}$ as follows:

$$DL_{similarity} = \{(s, sameAs, t) \mid s \in S \land t \in T \land d_{similarity}(s, t) < \theta_{d_{similarity}}\}$$





Link Discovery is the fourth and the most important Linked Data **Principle**.

Establish **semantic relations** between entities in order to **enrich the information** that is known about them. [Bizer et al., IJSWIS'06]







Let *S* and *T* be two sets of entities and *R* the set of relations that can be discovered between entities. For a relation $r \in R$, w.l.o.g., we define a distance function d_r and a distance threshold θ_{d_r} as follows:

$$d_r: S \times T \rightarrow [0,1]$$
, $\theta_{d_r} \in [0,1]$

We define the set of discovered links for relation r (DL_r) as follows:

$$DL_r = \{ (s, r, t) \mid s \in S \land t \in T \land d_r(s, t) < \theta_{d_r} \}$$



Link Discovery (Example)



Fields - OSM Water Bodies

Natura (2000) - Fields

LE®

01/06/2015

Discovering Spatial and Temporal Links among RDF Graphs



Link Discovery (Example)



Natura (2000) - Fields

Fields - OSM Water Bodies



Discovering Spatial and Temporal Links among RDF Graphs





- Different Data Providers create Heterogeneous
 Datasets
 - Example: Literal Heterogeneity (case, language, etc).

name = PORTOROZ

name = Portorose

- We focus on:
 - Heterogeneity in the Representation of Geospatial Information in RDF
 - Heterogeneity in the Representation of Temporal Information in RDF

Heterogeneity in the Representation of Geospatial Information in RDF

- _:1 rdf:type geo:Geometry .
- _:1 geo:hasGeometry

"<http://www.opengis.net/def/crs/EPSG/0/4326>

POINT(10 20)"^^geo:wktLiteral .

- _:1 rdf:type strdf:Geometry .
- _:1 strdf:hasGeometry

"<gml:Point crsName="EPSG:2100"><gml:coordinates>10,20
</gml:coordinates></gml:Point>"^^strdf:GML .

- _:1 rdf:type wgs84Geo:Point .
- _:1 wgs84Geo:lat "10"^^xsd:double .
- _:1 wgs84Geo:long "20"^^xsd:double .



Geosparn





Heterogeneity in the Representation of Geospatial Information in RDF

- _:1 rdf:type geo:Geometry .
- _:1 geo:hasGeometry

"<http://www.opengis.net/def/crs/EPSG/0/4326>

POINT(10 20)"^^geo:wktLiteral .

- _:1 rdf:type strdf:Geometry .
- _:1 strdf:hasGeometry

"<gml:Point crsName="EPSG:2100"><gml:coordinates>10,20

</gml:coordinates></gml:Point>"^^strdf:GML .

• Different Vocabularies









Heterogeneity in the Representation of Geospatial Information in RDF

- "<http://www.opengis.net/def/crs/EPSG/0/4326>
 POINT(10 20) ^^geo:wktLiteral .
 _:1 rdf:type strdf:Geometry .
 _:1 strdf:hasGeometry
 "<gml:Point crsName="EPSG:2100"><gml:coordinates>10,20
 </gml:coordinates></gml:Point>"^strdf:GML .
 - Different Vocabularies

:1 rdf:type geo:Geometry .

:1 geo:hasGeometry

• Different Serializations of Geometries











- Different Serializations of Geometries
- Geometries expressed in Different Coordinate Reference Systems (CRS)



Geospatial Information in RDF



LE

meldies

Geospatial Information in RDF



- Different Sampling Values
- Different Granularity
- Different Rounding Effects



Heterogeneity in the Representation of **Temporal Information in RDF**

:1 ex:hasBirthday "1989-09-24T11:05:00+01:00"xsd:dateTime

:1 ex:hasAffiliation ex:UoA "[2007-10-15T00:00:00+03:00,2013-10-15T00:00:00+04:00) "^^strdf:Period



melodies



Heterogeneity in the Representation of Temporal Information in RDF

_:1 ex:hasBirthday "1989-09-24T11:05:00+01:00"xsd:dateTime

_:1 ex:hasAffiliation ex:UoA

"[2007-10-15T00:00:00+03:00, 2013-10-15T00:00:00+04:00)"^^strdf:Period

• Different Vocabularies











- Different Vocabularies
- Different Time Zones



- Different Vocabularies
- Different Time Zones
- Time Instants and Periods





- Introduction to Entity Resolution and Link Discovery
 - Examples, Definitions, Common Problems
- Spatial Entity Resolution
- Spatial and Temporal Link Discovery
 - Background and Developed Methods
 - Extensions to the Silk Framework
 - Hands-on

Spatial Entity Resolution (Example Revisited)



Problem of understanding that two (or more) entities in **data-world** are references of the same **real-world** entity. [Christen, TKDE'11]



Discovering Spatial and Temporal Links among RDF Graphs

ISehgal et al. GIS'061

- Location Name Similarity
 Edit, Jaccard distance
- Location Similarity
 Euclidean distance
- Location Type Similarity
 - (e.g. type "river" is similar to type "stream")

Combines the above similarities to compute the overall similarity between entities





- Similarity measure: Hausdorff Distance
 - Intuitively Hausdorff Distance is defined as the largest distance between the closest points of two geometric shapes

$$l_{H}(A,B) = max \left\{ \begin{cases} \arg\max\arg\min_{a \in A} \arg\min_{b \in B} d(a,b), \arg\max\arg\min_{a \in A} d(a,b) \\ b \in B & a \in A \end{cases} \right\}$$



- Handling Geospatial Heterogeneity
 - Converts geometries to a common vocabulary (NeoGeo)
 - Assumes WGS-84 CRS
- Optimization
 - Simplifies Geometries with Ramer-Douglas-Peucker algorithm



Spatial Entity Resolution (3/4)



- Heuristic Combination of:
 - URI Similarity
 - Label Similarity
 - Considering the language of the labels
 - Location Similarity
 - Assuming the W3C Geo vocabulary
 - Geometric Similarity
 - Minimum Distance between two Geometries







- Non-Spatial Criteria
 Implemented within the LIMES framework
- Geometric Similarity
 - Hausdorff Distance
 - Optimizations
 - Bounding Circle: Avoids useless comparisons $\mu(s, t) = \delta(\zeta(s), \zeta(t)) - r(s) - r(t) > \theta \Rightarrow \delta(s, t) > \theta$
 - Space tiling: Reduces the quadratic number of comparisons





Spatial Entity Resolution



- [Sehgal et al. GIS'06]
 - Spatial and non-Spatial Criteria
 - Only Location Similarity
- [Salas et al., TerraCognita'11]
 - Only Spatial Criteria
 - Complex Geometric Similarity Methods
- [Vilches-Blázquez et al., AGILE'12]
 - Spatial and non-Spatial Criteria
 - Simple Geometric Similarity Methods
- [Ngonga Ngomo, ISWC'13]
 - Spatial and non-Spatial Criteria
 - Complex Geometric Similarity Methods
 - Reduced number of comparisons





- Introduction to Entity Resolution and Link Discovery
 - Examples, Definitions, Common Problems
- Spatial Entity Resolution
- Spatial and Temporal Link Discovery
 - Background and Developed Methods
 - Extensions to the Silk Framework
 - Hands-on





Link Discovery is the fourth and the most important Linked Data **Principle**.

Establish **semantic relations** between entities in order to **enrich the information** that is known about them. [Bizer et al., IJSWIS'06]







- Dimensionally Extended 9-Intersection Model [Clementini et al., SSD'93]
 - Captures topological relations in R², by considering the dimension (dim) of the intersections involving the interior (I), the boundary (B) and the exterior (E) of the two geometries.

$$\text{DE-9IM}(\mathbf{a},\mathbf{b}) = \begin{bmatrix} \dim(I(a) \cap I(b)) & \dim(I(a) \cap B(b)) & \dim(I(a) \cap E(b)) \\ \dim(B(a) \cap I(b)) & \dim(B(a) \cap B(b)) & \dim(B(a) \cap E(b)) \\ \dim(E(a) \cap I(b)) & \dim(E(a) \cap B(b)) & \dim(E(a) \cap E(b)) \end{bmatrix}$$

 Examples: Intersects, Equals, Touches, Disjoint, Contains, Crosses, Covers, CoveredBy and Within





- Region Connection Calculus [Randell et al. KR'92]
 - RCC-8: a well-known subset of RCC, which is based on eight topological relations



 DC stands for DisConnected, EC for Externally Connected, TPP for Tangential Proper Part, NTPP, for Non Tangential Proper Part, and TPPi and NTPPi are the inverse relations of TPP and NTPP





- Allen's Interval Calculus [Allen, Commun. ACM'83]
 - thirteen jointly exclusive and pairwise disjoint qualitative relations

Relation	Illustration
X before Y	X
Y after X	Y
X meets Y	X
Y isMetBy X	Y
X overlaps Y	X
Y isOverlappedBy X	Y
X starts Y	X
Y isStartedBy X	Y
X during Y	X
Y contains X	Y
X finishes Y	X
Y isFinishedBy X	Y
X equals Y	X Y





• We consider the previous Spatial (R_s) and Temporal (R_t) relations as Boolean relations (R_B) i.e., either they hold or they do not:

$$R_s, R_t \subset R_B$$

• R_B constitutes a special subset of R. The distance function d_r and the distance threshold θ_{d_r} for a relation $r \in R_B$ are defined as follows:

$$d_r(\mathbf{s},\mathbf{t}) = \begin{cases} 0 & if \ r \ holds \\ 1 & elsewhere \end{cases}, \ \theta_{d_r} = 1$$



Spatial and Temporal Transformations (1/2)



- **CRS Transformation**. The geometries of a dataset can be expressed in a Coordinate Reference System that is more precise for the geographic area that they describe (e.g., the GGRS87 for Greece). This transformation converts the CRS of a geometry to the World Geodetic System (WGS 84)
- Vocabulary Transformation. This transformation converts geometry literals from GeoSPARQL, stRDF or W3C GEO to a common vocabulary (GeoSPARQL)
- Serialization Transformation. This transformation converts the geometries of a dataset to a common serialization (WKT)
- **Time-Zone Transformation**. This transformation converts the time zone of a given time interval to Coordinated Universal Time (UTC)
- **Period Transformation**. This transformation converts a time instant to a period with the same starting and ending point



Spatial and Temporal Transformations (2/2)



- **Simplification Transformation**. Some datasets have very complex geometries, which makes the computation of spatial relations inefficient. This transformation simplifies a geometry according to a given distance tolerance, ensuring that the result is a valid geometry having the same dimension and number of components as the input
- Envelope Transformation. This transformation computes the envelope (i.e., the minimum bounding rectangle) of a geometry and it is useful in cases that we want to compute approximate spatial relations between two datasets
- Area Transformation. In some cases it is enough to compare just the areas of two geometries to infer whether they are the same or not. This transformation computes the area of a given geometry in square metres
- **Points-To-Centroid Transformation**. In crowdsourcing datasets like OpenStreetMap, multiple users can define the position of the same placemark. As a better approximation of the real position of this placemark we can compute the centroid of these positions. This transformation computes the centroid of a cluster of points





- Cartesian Product Technique (Naive)
 - Performs exhaustive checks between the pairs of the entities of datasets
 - Complete
 - Complexity: O(|S||T|) checks
- Blocking Technique [Isele et al., WebDB'11, Papadakis et al, TKDE'13]
 - Divides the entities into blocks
 - Decreases the number of checks
 - Complete
 - Complexity: O(|S||T|) checks (worst case), O(|L|) checks (best case)
- * |S|, |T|: number of entities in datasets S and T; |L|: number of links between datasets S and T





- Divide the surface of the earth into **curved rectangles** (blocks)
- Adjust the area of the blocks with a **blocking factor** (bf) (blockArea: $\frac{1}{bf^2}^{o^2}$)



- If the MBB of a geometry spatially intersects with a block, then insert it in this block
- Check for a spatial relation only within each block (independently)
- Construct the set of discovered links (DL_r) by **aggregating** the respective links that have been discovered within each block



- Divide the time into intervals (blocks)
- Adjust the length of the blocks with a **blocking factor** (bf) (blockLength: $\frac{1}{bf}$ time units)



- If a time period or instant temporally intersects with a block, then insert it in this block
- Check for a temporal relation only within each block (independently)
- Construct the set of discovered links (DL_r) by **aggregating** the respective links that have been discovered within each block



Blocking Technique







- Fully parallelizable with respect to the blocks
- Proven sound and complete
- 100% accurate links
- 100% precision, recall, F-measure



Extensions to the Silk Framework: Spatial and Temporal Relations





Discovering Spatial and Temporal Links among RDF Graphs



Buffer Serialization Area GeometryLiteral Points-To-Centroid Envelope Simplification TimeZone CRS



Extensions to the Silk Framework





- Spatial and Temporal Extensions for Silk implemented as Plugins
- **Transparent** to all the applications of Silk
 - Single Machine
 - MapReduce
 - Workbench





- Download: <u>https://github.com/silk-framework/silk</u>
- Workbench application pre-installed in the VM
- Discover the following links:

Source Dataset	Relation	Target Dataset	
Field Boundaries	Contains	Raster Cells	Easy
OSM Water Bodies	Intersects	Natura (2000)	Bonus
Natura (2000)	Within	Federal States of Germany	Bonus

All the datasets will be first converted to RDF with GeoTriples!





- [Bizer et al., IJSWIS'06]
 Bizer, C., Heath, T., Berners-Lee, T.: Linked Data The Story So Far. International Journal on Semantic Web and Information Systems 5(3), 1–22 (2009)
- [Christen, TKDE'11]

P. Christen, " A survey of indexing techniques for scalable record linkage and deduplication." in IEEE TKDE 2011.

• [Auer, RW'13]

Auer, S., Lehmann, J., Ngomo, A.C.N., Zaveri, A.: Introduction to Linked Data and Its Lifecycle on the Web. In: Rudolph, S., Gottlob, G., Horrocks, I., van Harmelen, F. (eds.) Reasoning Web. Lecture Notes in Computer Science, vol. 8067, pp. 1–90. Springer (2013)

• [Salas et al., TerraCognita'11]

Salas, J., Harth, A.: Finding spatial equivalences accross multiple RDF datasets. In: Proceedings of the Terra Cognita Workshop on Foundations, Technologies and Applications of the Geospatial Web. pp. 114–126. Citeseer (2011)

• [Sehgal et al. GIS'06]

Sehgal, V., Getoor, L., Viechnicki, P.D.: Entity resolution in geospatial data integration. In: Proceedings of the 14th annual ACM international symposium on Advances in geographic information systems. pp. 83–90. ACM (2006)



References (2/3)



• [Vilches-Blázquez et al., AGILE'12]

Vilches-Blázquez, L.M., Saquicela, V., Corcho, O.: Interlinking geospatial information in the web of data. In: Bridging the Geographic Information Sciences, pp. 119–139. Springer (2012)

 [Ngonga Ngomo, ISWC'13] Ngonga Ngomo, A.C.: Orchid - reduction-ratio-optimal computation of geo-spatial distances for link discovery. In: Proceedings of ISWC 2013 (2013)

• [Clementini et al., SSD'93]

Clementini, E., Di Felice, P., van Oosterom, P.: A small set of formal topological relationships suitable for end-user interaction. In: Abel, D., Chin Ooi, B. (eds.) Advances in Spatial Databases, Lecture Notes in Computer Science, vol. 692, pp. 277–295. Springer Berlin Heidelberg (1993), http://dx.doi.org/10.1007/3-540-56869-7_16

- [Randell et al. KR'92] Randell, D.A., Cui, Z., Cohn, A.G.: A spatial logic based on regions and connection. In: KR. pp. 165–176 (1992)
- [Allen, Commun. ACM'83] Allen, J.F.: Maintaining knowledge about temporal intervals. Commun. ACM 26(11), 832– 843 (Nov 1983)



References (3/3)



- [Isele et al., WebDB'11]
 Isele, R., Jentzsch, A., Bizer, C.: Efficient multidimensional blocking for link discovery without losing recall. In: WebDB. Citeseer (2011)
- [Papadakis et al, TKDE'13] Papadakis, G., Ioannou, E., Palpanas, T., Niederée, C., Nejdl, W.: A blocking framework for entity resolution in highly heterogeneous information spaces. Knowledge and Data Engineering, IEEE Transactions on 25(12), 2665–2682 (2013)





Thanks for your attention! Questions?