Knowl. Org. 46(2019)No.8

607

Shu-Jiun Chen. Semantic Enrichment of Linked Personal Authority Data: A Case Study of Elites in Late Imperial China

# Semantic Enrichment of Linked Personal Authority Data: A Case Study of Elites in Late Imperial China†

## Shu-Jiun Chen

Institute of History and Philology, Academia Sinica, Taipei 11529, Taiwan,
<sophy@sinica.edu.tw>

**Abstract:** The study uses the Database of Names and Biographies (DNB) as an example to explore how in the transformation of original data into linked data, semantic enrichment can enhance engagement in digital humanities. In the preliminary results, we have defined instance-based and schema-based categories of semantic enrichment. In the instance-based category, in which enrichment occurs by enhancing the content of entities, we further determined three types, including: 1) enriching the entities by linking to diverse external resources in order to provide additional data of multiple perspectives; 2) enriching the entities with missing data, which is needed to satisfy the semantic queries; and, 3) providing the entities with access to an extended knowledge base. In the schema-based categories that enrichment occurs by enhancing the relations between the properties, we have identified two types, including: 1) enriching the properties by defining the hierarchical relations between properties; and, 2) specifying properties' domain and range for data reasoning. In addition, the study implements the LOD dataset in a digital humanities platform to demonstrate how instances and entities can be applied in the full texts where the relationship between entities are highlighted in order to bring scholars more semantic details of the texts.

## 1.0 Introduction

Semantic relations between entities are the basic units of knowledge organization (Green 2001; Stock 2010). This paper discusses the issue of semantic enrichment in linked open data (LOD) research. The study will use the Database of Names and Biographies (DNB) as an example to explore how, in the transformation of original data into linked data, semantic enrichment can facilitate research inquiries and enhance engagement in digital humanities. Hosted by the Institute of History and Philology, Academia Sinica (Taiwan), the DNB contains 35,666 records of Chinese historical persons who are cultural and sociopolitical elites in late imperial China (1368-1911), extracted from various historical archives for the purpose of sup-porting historians' research. Each metadata record has information including name, alternative name, dates of birth/death, native place, biographical data, work experiences, related persons, specialties, academic background, job titles and references. Semantic enrichment is one of the current LOD project's core research streams, developed by collaborating with historians for converting the DNB from legacy databases to linked data format for the purpose of digital humanities. The study has deployed the methods of data modeling, data reconciliation and data enrichment during transforming the legacy metadata records into LOD in order to add value that is structured in a machine-processable format and gives more meaning to the dataset (Hyvönen 2018; Van Hooland and Verborgh 2014; Zeng 2019). The DNB data model is composed of five

core classes (i.e., agent, event, place, object and time classes), sixty-seven properties from sixteen semantic vocabularies (i.e., bio, dbpedia, dcterms, gvp, leo, owl, rdfs, schema.org, skos, etc.) and reuses three external resources (i.e., AAT, VIAF, TGAZ). The LOD dataset contains more than two million triples and is available to query linked data with SPARQL from the LODLab of Academia Sinica (http://data.ascdc.tw/en/sparql.php). The following presents the preliminary results of the study, which we have defined instance-based and schema-based categories of semantic enrichment.

## 2.0 Instance-based semantic enrichment

Instance-based, or entity, enrichment, in the context of LOD is a process to endow the original data with other supplementary or supported knowledge, which is beyond the original data content, from external resources. This is done to extend the perspective of the data itself and also integrate heterogeneous resources into a more complete and sophisticated knowledge. By enriching the entities in the LOD-based dataset, it not only enlarges the knowledge base itself but also inspires new perspectives for further research and interpretation of the study results. There are three approaches identified in the study to achieve entity enrichment as follows.

### 2.1 Enriching the entities by linking to diverse, cross-domain external resources in order to provide additional data of multiple perspectives

In the DNB dataset, which is focused on the information of Chinese historical figures, the agent entity is linked to the cross-domain external resources. For instance, in the agent entity of Tseng Guo-fan (曾國藩, 1811-1872), a famous literati and high minister of the late Qing period, the

AAT (Art & Architecture Thesaurus) concept of "calligraphy" is directly linked with the agent by the specialty property (dbpedia-owl:specialty). In the original DNB metadata, a specialty of Tseng is presented in literal as"書法" (calligraphy) in Chinese characters. Since the AAT is a multilingual thesaurus including English, Chinese, German, Dutch and more (Harpring 2018), the link to the AAT term for "calligraphy" (AAT 300053162) can enrich the agent entity in the related property with information on the same concept and its definition in other languages. This could enhance understandability for users who do not know the Chinese language.

### 2.2 Enriching the entities with missing data which is needed to satisfy the semantic queries

As to the information on the native place of a person, since its value shows the historical administrative name of a place and is not entirely equivalent with the current place name, we, therefore, added a new contextual entity for a place for each agent when the referred agent has information in the field of "native place." The method of adding a place entity is not only to describe the historical name but also to be enriched by linking to the linked data of the Chinese historical place name in the Temporal Gazetteer (TGAZ) developed by Harvard University, which defines the geographical range and temporal information of the related historical place name. Such a place entity is linked to the agent entity of this study by the native-place property (ascdc:nativePlace) and also carries the label (rdfs:label) and the link to external resources by the same-as property (owl:sameAs) (shown in figure 1).

Another reason to reuse the terms of Chinese historical place names from TGAZ is that on the SPARQL (SPARQL Protocol and RDF Query Language) interface for DNB dataset, the study applied Chinese historical Geographic Information System (GIS) maps for showing an agent's native
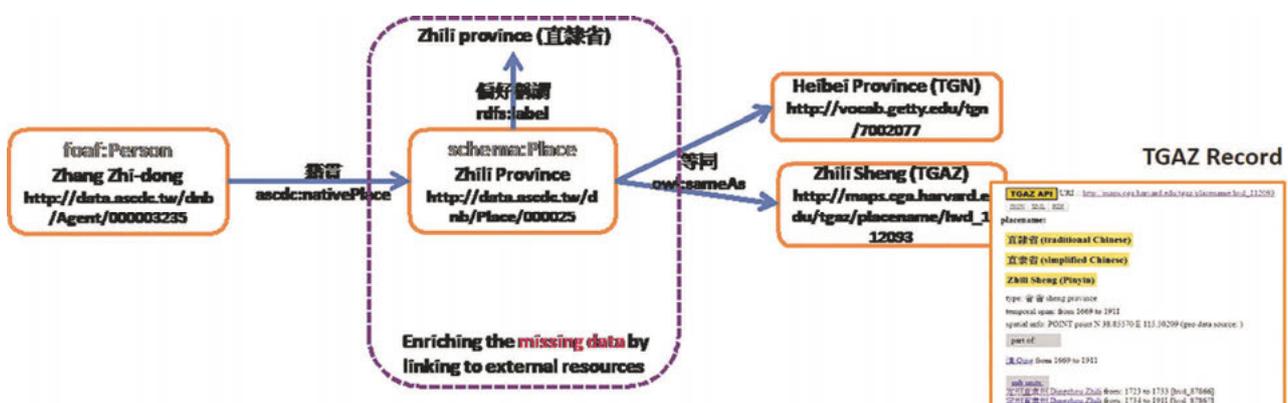


*Figure 1.* Enriching the entities with missing data, which is needed to satisfy the semantic queries.

609

Knowl. Org. 46(2019)No.8
Shu-Jiun Chen. Semantic Enrichment of Linked Personal Authority Data: A Case Study of Elites in Late Imperial China

place in certain examples of SPARQL query. For instance, when one queries "which provenances did the Qing agents ranked as jinshi (進士, "presented scholar," the highest degree of the Chinese Imperial civil service examinations) come from?" Since the current Google Map as open source shows only the global map in modern political and administrative boundaries, it is quite different from the historical map. The location of a city from hundreds of years ago might also be different from its current geographic situation. Since each term in the Temporal Gazetteer contains the longitude and latitude information of a historical place name in different dynasties or periods, the application of TGAZ is reasonable for displaying the related geographic information in the historical GIS-maps.

## 2.3 Providing the entities with extension to a knowledge base

The extension of data content in the entities to a knowledge base is regarded in this study as one part of data enrichment. It could not only enable the entities with more detailed information from the external resources but also make linkage of the original entity to an entire knowledge base of a certain subject field, which can inspire new perspectives for further research and interpretation of the study results. For instance, the study makes extension of entities to the time ontology, which is one part of the knowledge base. In the current LOD-based dataset of the DNB, the entities relating to the temporal information on the official career of a person is particularly extended to the "term lists of the time and periods," a con-

trolled vocabulary developed by the study for defining the Chinese historical periods and yearly times of all dynasties in China (see Figure 2). The application of such term lists is based on the "time ontology," which is developed by the study and based on W3C's time ontology in OWL (Cox and Little 2017).

For example, the DNB person agent Ding Bao-zhen (丁寶楨) was appointed as Governor of the Sichuan Province (四川總督) between 1881 and 1886. In the data model design of DNB, a person agent's official career in the government is expressed as entities of official service (ascdc:OfficialService). To describe the beginning date and end date of Ding's appointment, this entity for official service is linked to an entity of temporal period (time:TemporalEntity) by at-time property (leo:atTime), which further describes the beginning and end date of the related period for official service in the historical Chinese era year, represented by the "time:instant" entity. To extend the knowledge base for that temporal information, the study's time ontology is applied to describe the hierarchical temporal details of the instances as the "7th year of Guangxu" (光緒七年, 1881) and "12th year of Guangxu" (光緒12年, 1886) and linked with the entity for era name (光緒), emperor name (光緒皇帝) and dynasty name (清朝) in "time:propertInterval" by the properties "time:intervalDuring" and "time:inside"(see Figure 3).

After extending the entity with the temporal terms to describe the beginning and ending year of a certain period of an official position, the information of these mentioned historical years is expressed as entities of thing and can further be linked to the era names, emperor names and
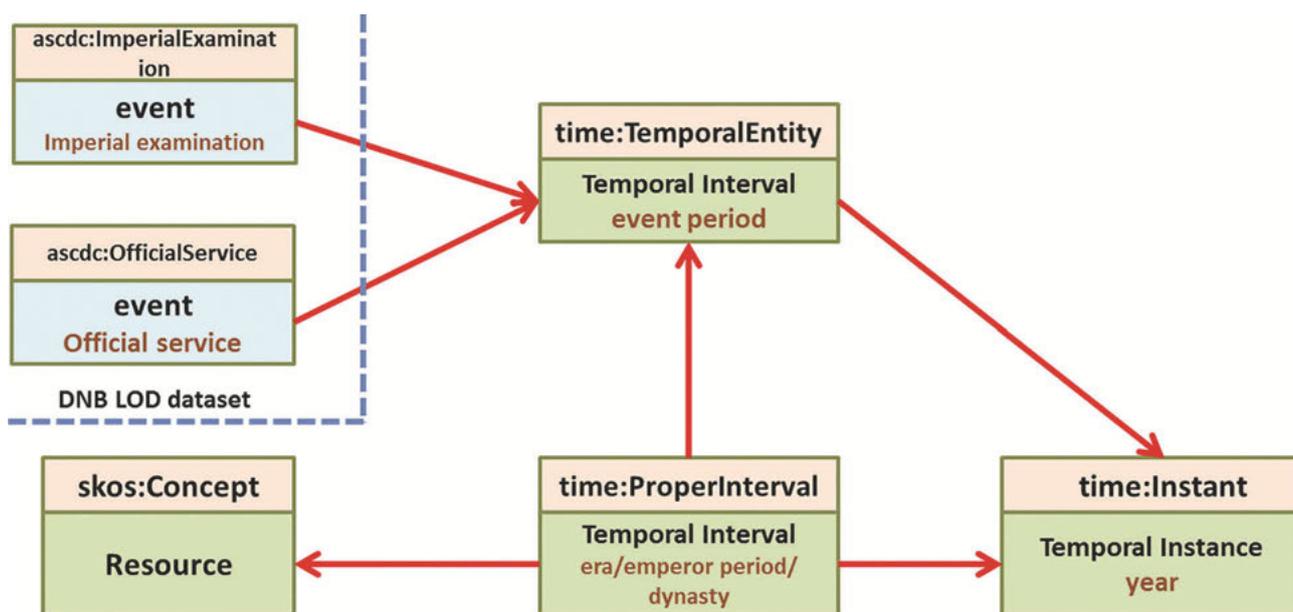


*Figure 2.* Extension of DNB entities with the study's time ontology, viewed with classes.

610

Knowl. Org. 46(2019)No.8
Shu-Jiun Chen. Semantic Enrichment of Linked Personal Authority Data: A Case Study of Elites in Late Imperial China

the dynasty names in Chinese history. Therefore, a hierarchically structured knowledge base, which is focused on Chinese historical temporal names, is entirely integrated into the dataset and enriches the data content of each related person agent.

### 3.0 Schema-based semantic enrichment

Schema- or property-based enrichment in LOD is a process to enable hierarchical or associative meaning within pairs of property, which could create a relationship between related entities and also enable a meaningful and efficient data query in a hierarchical semantic structure.

### 3.1 Enriching the properties by defining the hierarchical and associative relations between properties

To a certain extent, the applicability of enriching the properties in the LOD datasets depends on whether a hierarchical or associative meaning exists between different data elements of the original metadata. In the data element of the current DNB datasets, such related meaning is especially found in the data element as the "personal relations" (人物關係) of an agent, in which a person's connections to another agent is linked and expressed by reusing suitable properties. As an example, the semantic relationships of an agent A to his teacher, agent B, is defined as the has-teacher property (agrelon:hasTeacher), while the relationships of agent A to his grandparent, agent C, is expressed as the has-grandparent property (agrelon:hasGrandparent).

In fact, the personal relation between agents is a mutual- or hierarchical-expressible relationship. If A is student of B, then the B should be the teacher of A. In the semantic data model design, such mutual or hierarchical relation can be defined by enriching the definition of the related properties. In the current world of semantic web, RDFS and OWL are the two types of data vocabularies, which are mostly applied to enrich the relations between properties as the abovementioned cases and also to enhance the efficiency by data reasoning (Polleres, Axel, et al. 2013). In particular, the subproperty-of property (rdfs:subpopertyOf) can be used to mark the hierarchical relation between properties, while the inverse-of property (owl:inverseOf) is suitable to describe the mutually affected relations between properties on the same level (see Table 1).

Taking the agent Zeng Guo-fan (曾國藩, DNB: NO000000058) in the DNB datasets as an example, the figure is linked to the agent Li Hong-zhang (李鴻章, DNB: NO000002242) by the has-student property ("agrelon:hasStudent). Since the relationship between a teacher and student is mutually referred, if the inverse-of property (owl:inverseOf) is reused and enriches the definition of the has-student property (agrelon:hasStudent), the reverse relation expressed as has-teacher property (agrelon:hasTeacher) would also be findable by data reasoning.



*Figure 3*. Extension of the entity with knowledge-based external resource (Example: Extension of the entity for official service with hierarchical information on the beginning and end date of Ding Bao-zhen's appointment as Governor in Sichuan between 1881 and 1886).

| | Properties | Domain | Range | Function of property enriching |
|---|---|---|---|---|
| 1 | **rdfs:subPropertyOf** | Property | Property | Hierarchical relation |
| 2 | **owl:inverseOf** | Property | Property | Mutual relation |

*Table 1*. Enrich the properties by using rdfs:subpopertyOf and owl:inverseOf.

With the same example of Zeng Guo-fan, the figure is further linked to the agent Yuan Bingzhen (袁秉楨, DNB: NO000012514) by the has-child-in-law property (agrelon:hasChildInLaw) and to Zeng Guangquan (曾廣銓, DNB: NO000008193) by the has-grandchild property (agrelon:hasGrandchild). Since those different types of relations to a person can be clustered in a broader range of properties such as relatives, a hierarchical structure of property can be hence defined by using the sub-property-of property (rdfs:subPropertyOf) to structure the has-child-in-law property (agrelon:hasChildInLaw) and the has-grandchild property (agrelon:hasGrandchild) both under the property as has-relative (agrelon:hasRelative).

### 3.2 Specifying properties' domain and range for data reasoning

In the semantic web, each data can be expressed as a triple, which is composed of subject, property and object. From them, the major function of a property is to enable the entities of information (subject and object) with a semantic relation, which could enable the data query in a logical, machine-processable and machine-understandable way. In other words, the property plays a role as the bridge to connect the subject with object and thus construct complete, meaningful information in the data. However, the use of a certain property is not arbitrary. Each property has its own definition to which condition or restriction can be applied to link the subject- and object-entities.

In the current data model design for the "Database of Names and Biographies" (DNB), sixty-seven properties from sixteen vocabularies are reused to describe the biographic information on Chinese historical figures and relations between figures. The information on the domain and range of all reused properties is described in the specification of DNB ontology, seen in the selected examples in Table 2. Such specification can be used as referential

standard by the semantic data model design and also defines which data context that a property could be reused in the data structure.

In the semantic data model design, the domain of a property is always the instance of an entity (or a class), which defines the subject of an information. As in the abovementioned example, a property as "bio:father" in the DNB is applied to describe the information on the father of a person. The domain of this property is defined as "foaf:Agent," which means this property is only suitable to be applied by an entity of the agent. In the machine-processable form, it could be formulated as "rdfs:domain foaf:Agent." However, the range of a property could be expressed as a different data type, such as the instance of an entity (or a class), literal information, date decimal numbers of measurement or quantity of item. Again, as in the example (Table 2), the specification also defines the range of "bio:father" as "foaf:Agent," which means that object information should also be an instance of agent entity. In the machine-processable form, it could be formulated as "rdfs:range foaf:Agent."

### 4.0 Applications of LOD for digital humanities

The input of the LOD-based dataset of the Database of the Names and Biographies (DNB) into the system of the Digital Humanities Research Platform (DHRP), developed by the Academia Sinica Center for Digital Cultures (ASCDC), is a linked data application to enhance the reusability of the DNB data. Additionally, the study demonstrates the possibility of applying a LOD-based dataset to enlarge the research scope of the scholars in digital humanities and to integrate into digital research tools- using examples in the DHRP.

The DHRP is an open, cloud-based text repository to enhance the research of digital humanities, which is developed as a platform for online services based on the needs of scholars. The platform is equipped with different digital tools for text and visual analytics, such as text annotation,

| bio:father | |
|---|---|
| URI | http://purl.org/vocab/bio/0.1/father |
| Label | Father |
| Type | Property |
| Comment | To describe information on father of a DNB Agent Entity |
| Domain | foaf:Agent |
| Range | foaf:Agent |
| Quantification | 0-1 |
| Data type | Concept/ASCDC |
| Examples | Liu yun father Liu Ton-hsung |

*Table 2*. Specifying domain and range of the father property (bio:father) in DNB.

text similarity comparison, N-gram analysis, historical spatiotemporal visualization or social network analysis. The digital content of the DHRP is currently uploaded with texts from rare Chinese books for a total of more than 220 million words. In the current stage, the DNB dataset is already in the test version of the DHRP-platform. In particular, we use the DNB's data on properties of the person's relations, specialty and native place as practical cases of studies to map the text passages in the *Qing Shilu* (the Veritable Records of the Qing Dynasty/清實錄) and to demonstrate the semantic relations between different person agents or agent entities with place or concept entities in the historiographical works of the DHRP-platform. In total, more than 20,000 named-entities from the 93,431 text passages in the *Qing Shilu* are matched with the instances of entities in the DNB.

In the DHRP-platform, the mapped DNB personal names in the *Qing Shilu* will be marked up in different colors according to their types in the data unit of a triple in the DNB dataset For instance, the person's name belonging to the subject in a triple will be highlighted in blue, while names of an object in a triple will be shown upon a gray background (Figure 4). The type of semantic relationships (properties) between different agents will be represented in a dotted line, which link the subject entity with its related object entity.

When moving the mouse cursor onto the subject agent in the text, the platform will automatically present the related agents of person names in green with the type of semantic relation. Further, clicking on those names will direct to the website of the LOD-datasets, showing the data content of the related records of the persons (Figure 5).

For further presentation of the related data in DHRP-platform by using the tools for data visualization, the function of social network analysis (SNA) is integrated into the system to show the matched persons in *Qing Shilu* based on the DNB dataset (Figure 6).

In the original DHRP-platform, the scholars could only execute the text retrieval and comparison based on the literal context uploaded in the system. Users could not find out detailed information or definitions of the retrieved words or text, since the context was not linked to the external resources by a sematic method. After uploading the LOD-based dataset in the originally text-based DHRP-
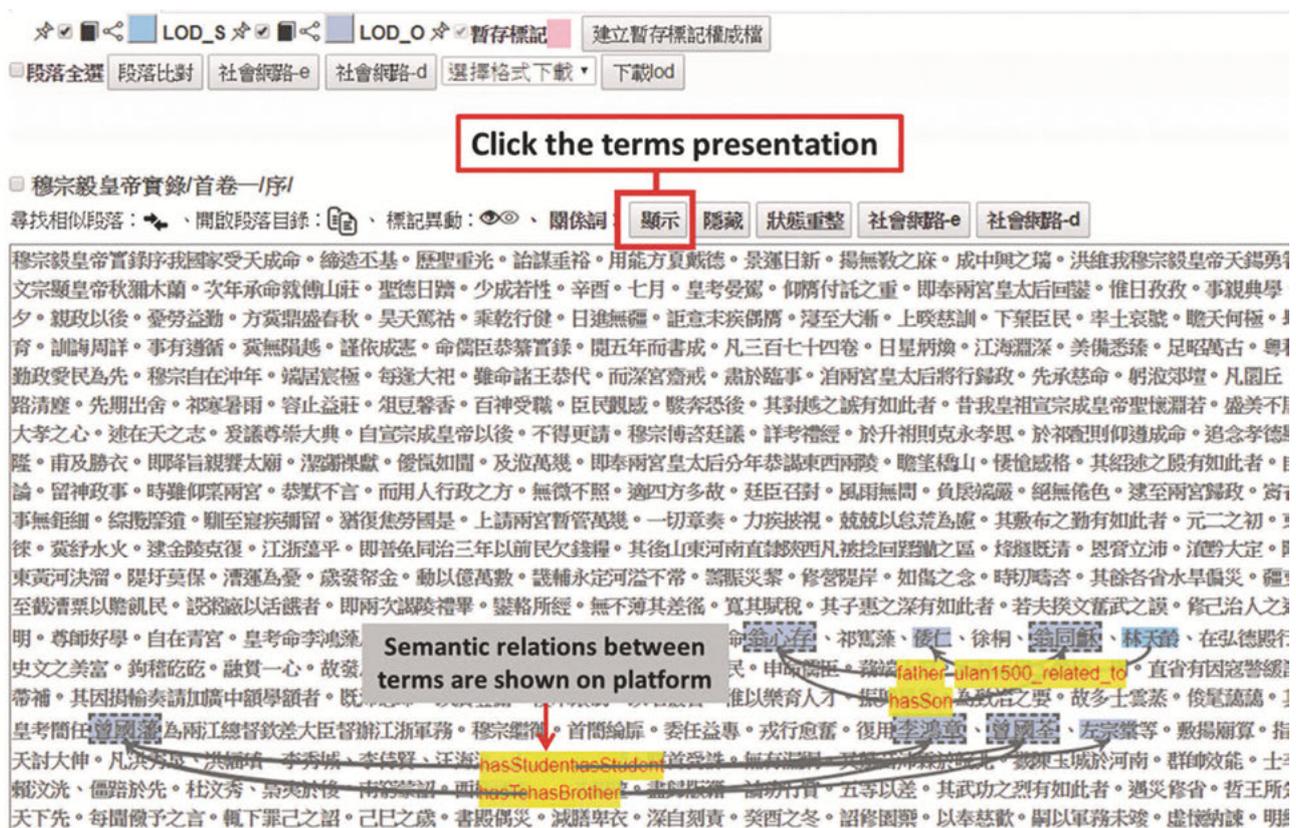


*Figure 4*. Revealing the relations between a matched person's name from DNB in *Qing Shilu* on the DHRP platform (Example: Teacher-student Relationship between Zeng Guo-fan (曾國藩) andd Li Hong-zhang (李鴻章) are shown with a type of sematic relation (has-Student) in red).

Knowl. Org. 46(2019)No.8
Shu-Jiun Chen. Semantic Enrichment of Linked Personal Authority Data: A Case Study of Elites in Late Imperial China

613



*Figure 5.* Linking the matched person's name in *Qing Shilu* to the external resource in the DNB (Example: Moving the cursor onto the person name of Zeng Guo-fan (曾國藩) and showing his related persons retrieved in *Qing Shilu*; clicking one of the names as Zeng Guo-chuan (曾國荃) and linking to the equivalent resource in the DNB).



*Figure 6.* SNA-analysis showing the relation types between matched persons in a passage of *Qing Shilu* in different forms of data visualization. (1) SNA by e-chart; (2) SNA by D3.js data visualization.

614

Knowl. Org. 46(2019)No.8

Shu-Jiun Chen. Semantic Enrichment of Linked Personal Authority Data: A Case Study of Elites in Late Imperial China

platform, a detailed definition of the matched text (for example, the further biographical information of a person) can be shown in DHRP by reusing the related data in DNB dataset. This is accomplished through the named-entity recognition, which is a mapping procedure of the terms in DNB with the text in the platform. The semantic relations of a person to another person, place or concept will also be notified by endowing the type of relations (properties) in the platform. These could extend not only the knowledge base of a scholar but might also offer other relevant information or inspire research angles, which one might not take notice when retrieving the results merely in a literal context.

## References:

Baca, Murtha and Melissa Gill. 2015. "Encoding Multilingual Knowledge Systems in the Digital Age: The Getty Vocabularies." *Knowledge Organization* 42: 232-43.

Bischof, Stefan. 2017. "Complementary Methods for the Enrichment of Linked Data." PhD diss., Technische Universität Wien. https://aic.ai.wu.ac.at/~polleres/supervised_theses/Stefan_Bischof_Dissertation_2017.pdf

Cox, Simon and Chris Little. 2017. "*Time Ontology in OWL.*" Retrieved from https://www.w3.org/TR/owl-time/

Ding, Li, Joshua Shinavier, Zhenning Shangguan and Deborah L. McGuinness. 2010. "SameAs Networks and Beyond: Analyzing Deployment Status and Implications of OWL:sameAs in Linked Data." In *The Semantic Web: ISWC 2010*, ed. Peter F. Patel-Schneider, Yue Pan, Pascal Hitzler, Peter Mika, Lei Zhang, Jeff Z. Pan, Ian Horrocks and Birte Glimm. Lecture Notes in Computer Science 6496. Berlin: Springer, 145-60.

Green, Rebecca. 2001. "Relationships in the Organization of Knowledge: An Overview." In *Relationships in the Organization of Knowledge*, ed. Carol A. Bean and Rebecca Green. Dordrecht: Springer, 3-18.

Harpring, Patricia. 2018. "Linking the Getty Vocabularies: The Content Perspective, Including an Update on CONA." In *Human Rights in Cyberspace: Proceedings of the 2018 Pacific Neighborhood Consortium Annual Conference and Joint Meetings,* ed. Shih-Lung Shaw, Ta-Chien Chan and Ling-Jyh Chen. Piscataway, NJ: IEEE. doi:10.23919/PNC.2018.8579460

Hyvönen, Eero. 2018. "Cultural Heritage Linked Data on the Semantic Web: Three Case Studies Using the Sampo Model." In *Datu ireki estekatuak eta informazioaren kudeaketa integrala ondare-erakundeetan: VIII. Arte Garaikideko Dokumentazio Zentroen Topaketak*. Vitoria-Gasteiz, Spain: Artium, 19-20.

Manguinhas, Hugo, ed. 2016. "Europeana Semantic Enrichment Framework: Documentation," Contributors: Antoine Isaac, Valentine Charles, Yorgos Mamakis, Juliane Stiller. Version: 17th November 2016 (updated 2017, 2018). https://docs.google.com/document/d/1JvjrWMTpMIH7WnuieNqcT0zpJAXUP06x4uMBj1pEx0Y/edit#

Polleres, Axel, Aidan Hogan, Renaud Delbru and Jürgen Umbrich. 2013. "RDFS and OWL Reasoning for Linked Data." In *Reasoning Web 2013: Reasoning Web. Semantic Technologies for Intelligent Data Access*, ed. Sebastian Rudolph, Georg Gottlob, Horrocks Ian and Frank van Harmelen. Lecture Notes in Computer Science 8067. Berlin: Springer, 91-149.

Stock, Wolfgang G. 2010. "Concepts and Semantic Relations in Information Science." *Journal of the American Society for Information Science and Technology* 61: 1951-69.

Subhashree, S., Rajeev Irny and P. Sreenivasa Kumar. 2018. "Review of Approaches for Linked Data Ontology Enrichment." In *Distributed Computing and Internet Technology: 14th International Conference, ICDCIT 2018, Bhubaneswar, India, January 11–13, 2018, Proceedings*, ed. Atul Negi, Raj Bhatnagar and Laxmi Parida. Lecture Notes in Computer Science 10722. Cham: Springer, 27-49.

Van Hooland, Seth and Ruben Verborgh. 2014. *Linked Data for Libraries, Archives and Museums: How to Clean, Link and Publish Your Metadata*. Chicago: Neal-Shuman.

Zeng, Marcia Lei. 2019. "Semantic Enrichment for Enhancing LAM Data and Supporting Digital Humanities. Review article." *El profesional de la información*, 28, no. 1. doi:10.3145/epi.2019.ene.03