Knowledge Graph extraction from Scholarly Data

Construction de graphes de connaissances à partir des publications scientifiques

Davide Buscaldi LIPN - Laboratoire d'Informatique de Paris Nord Université Paris 13

davide.buscaldi@lipn.univ-paris13.fr

Journée commune AFIA-ARIA, Paris 2/12/2019







Université Sorbonne

Plan of the talk

- Motivation

- Conclusion



- Proposed architecture
- Experiments and evaluation

•How difficult it is to ... compile a state of the art about a given topic? assess the novelty of a work? •find experts to review some papers?







- The production of scientific literature is growing at an increased pace
 - More conferences, venues, journals, etc.
 - New publishing paradigms that accelerate the process (ex: Open Access)
- Also: new research topics are continuously created

PubMed central growth, 2000-2013 (Ware and Mabe, The STM report 2015)



1980--2019

1980-1989

Title Unigram

Row	Paper-Title-Unigra	m
1	neural	1,390
2	task	1,070
3	learning	988
4	language	984
5	word	836
6	translation	828
7	machine	741
8	semantic	596
9	corpus	586
10	text	583
11	semeval	583
12	analysis	499
13	model	462
14	classification	448
15	models	427
16	embeddings	425
17	multi	419
18	networks	416
19	data	411
20	sentiment	408
		ОК 1К 2
		#papers =



Elaboration by Saif M. Mohammad on ACL Anthology data





Some (old) problems

Researchers may use the same word to refer to different topics

Or they may use different words to refer to the same topic

• Ex: different research domain • *lattice* in Physics vs. Computer Science • Ex: different research communities within the same domain • queries in relational databases vs. IR systems

• Ex: reaching the same result in parallel • Meucci's "teletrofono" vs. G.Bell's telephone



C \cap



Google Scholar

toponym disambiguation in information retrieval

Articles

Place name

Resolution

Date indifférente

Depuis 2018 Depuis 2017 Depuis 2014 Période spécifique...

Trier par pertinence

Trier par date

Toutes les langues

Rechercher les pages en Français

✓ inclure les brevets ✓ inclure les citations

Créer l'alerte

Conseil : Recherchez des résultats uniquement en Français. Vous pouvez indiquer votre langue de recherche sur hesis

Toponym disambiguation in information retrieval

D Buscaldi - 2010 - riunet.upv.es

In recent years, geography has acquired a great importance in the context of Information Retrieval (IR) and, in general, of the automated processing of information in text. Mobile devices that are able to surf the web and at the same time inform about their position are ...

 $\cancel{2}$ 99 Cité 12 fois Autres articles Les 17 versions \gg

Geographical information retrieval

CB Jones, RS Purves - Encyclopedia of Database Systems, 2009 - Springer ... Synonyms Place names; Toponyms; Knowledge organization sys- tems; Ontologies ... PostGIS Geographic Information System relies upon the PostgreSQL GiST implementation for its spatial ... GiST-based indexes have been used for applications in image retrieval, astronomy data ... $\cancel{2}$ 99 Cité 198 fois Autres articles Les 12 versions

Approaches to **disambiguating toponyms**

D Buscaldi - Sigspatial Special, 2011 - dl.acm.org

... In [6], the objective was to disambiguate toponyms in a local Italian newspa- per, where the granularity of toponyms was intended to be at the level of street ... However, there is still room for further work: toponym disambiguation may be used to tag ambiguous toponyms in Web ... $\cancel{2}$ 99 Cité 46 fois Autres articles Les 14 versions

A conceptual density-based approach for the **disambiguation** of **toponyms** D Buscaldi, P Rosso - ... Journal of Geographical Information Science, 2008 - Taylor & Francis ... 7. Conclusions and future work **Toponym disambiguation** is an open problem in GIR ... The obtained results expose the limits of both WordNet as a resource for the **disambiguation** of **toponyms** and of GeoSemCor as a resource for the testing of **disambiguation** methods ... $\cancel{2}$ $\cancel{2}$ Cité 89 fois Autres articles Les 5 versions



etrieval

Q

it en Français. Vous pouvez indiquer votre langue de recherche sur

ase Systems, 2009 - Springer

gazetteers must support multiple names and, ideally ... perspective into how regulation takes **place**, in ebi.ac.uk/ microarray-as/aer/ GEO http://www ...

2 versions

egions for spatial information retrieval

ce on Spatial Information ..., 2003 - Springer cimations of place name regions ... coordinates, there o-referenced through **place names** ... important aspects tic Web are the **resolution** of indi ...

versions

and spatial browsing

n [papers presented at ..., 1996 - ideals.illinois.edu aphic Information Re- trieval (GIR) in the context ... idexing and **retrieval** methods appropriate ... commonly Head- ings and **Name** Authorities as ...

3 versions ⇒>>

Name resolution in a directory database

Q





Not only Information Retrieval

Finding the scientific works related to a given topic is just a part of the problem

- •Which works are the most important?
- •Which works are the most innovative?
- •Who are the leading experts for a given topic?
- •What are the emerging topics in a research domain?
- Which methods have been applied to a particular task?

Dedicated search engines currently offer only a partial solution



Scientific Knowledge Graphs

Knowledge Base





Mining textual contents of research papers for build a SKG that represent the underlying knowledge

Data generation is no longer the limiting factor in advancing biological research. In addition, data integration, analysis, and interpretation have become key bottlenecks and challenges that biologists conducting genomic research face daily. To enable biologists to derive testable hypotheses from the increasing amount of genomic data, we have developed the VirtualPlant software platform. VirtualPlant enables scientists to visualize, integrate, and analyze genomic data from a systems biology perspective. VirtualPlant integrates genome-wide data concerning the known and predicted relationships among genes, proteins, and molecules, as well as genome-scale experimental measurements. VirtualPlant also provides visualization techniques that render multivariate information in visual formats that facilitate the extraction of biological concepts. Importantly, VirtualPlant helps biologists who are not trained in computer ...

Scientific Knowledge Graphs: an example







Our proposed architecture

- Formally, given a set of $D = \{d_1, \ldots, d_n\}$ scientific documents, we build a model $\gamma : D \rightarrow T$, where T is a set of triples (s, p, o)
 - s and o belong to a set of entities E and p belongs to a set of relation labels L
- Our framework includes the following steps:
 - 1. Extraction of entities and triples, combining various tools
 - 2. Entity refining, in which the resulting entities are merged
 - 3. *Triple refining,* in which the triples extracted by the different tools are merged together and the relations are mapped to a common vocabulary
 - 4. *Triple selection,* in which we select the set of «trusted» triples that will be included in the final output





Scientific Texts



Mining Scholarly Data for Scientific Knowledge Graph Construction







Examples of entities:

Semantic Web knowledge acquisition method theoretical learning **ODESeW**

Examples of extracted relations:

ODESeW, develop, ontology-base portals (OpenIE) Semantic e-Learning, enable, intelligent operations (OpenIE) theoretical learning, used-for, learning processes (Deep Learning Extractor) *machine readable information, part-of, Semantic Web* (Deep Learning Extractor)

- Six types of entities (Task, Method, Metric, Material, General Scientific, Other)
- Seven types of relations (Compare, Part-of, Conjunction, Evaluate-for, \bullet Feature-of, Used-for, Hyponym-Of)
- Problem: precise but the coverage is limited

- OpenIE (Stanford Core NLP) extracts open-domain relationships between any noun phrase
- CSO Classifier is a keyword extractor for classifying research papers in the CS domain

Mining Scholarly Data for Scientific Knowledge Graph Construction

Entity + Triple Extraction

Extractor Framework: neural model [Luan Yi et al.] based on SemEval-2018 task 7 winner system

Compensated by including more extractors into the process:

- CSO is paired with the Stanford NLP POS tagger to extract verbs between two keywords extracted by CSO Classifier





















(she; took; midnight train)

Born in a town, she took the midnight train

https://nlp.stanford.edu/software/openie.html

. . .

Mining Scholarly Data for Scientific Knowledge Graph Construction

OpenIE







CS Ontology

https://cso.kmi.open.ac.uk/home

computer science 👝

information retrieval o









https://cso.kmi.open.ac.uk/classify/

Mining Scholarly Data for Scientific Knowledge Graph Construction

CSO Classifier









Example of entity merging:

knowledge acquisition method knowledge acquisition approach

knowledge acquisition method

Mining Scholarly Data for Scientific Knowledge Graph Construction

Entity Refining and Merging

• Two or more entities may refer to the same concept for various reasons

• Refining:

- Acronym detection, punctuation and spaces, lemmatization...
- Removing "generic" entities based on domain frequency
- Splitting long entities (if "and" present

 Merge similar meanings (word embeddings) and synonyms (CSO) Ontology)

- If the cosine similarity is over a given threshold (0.85) the two
- entities are mapped on the same one







semantic web, <mark>analyze</mark> , markup language	analyze	>	w2	
semantic web, <mark>extend</mark> , markup language	extend		w3	
semantic web, <mark>combine</mark> , markup language	combine	>	w1	
semantic web, <mark>use</mark> , markup language	use	>	w4	
semantic web, <mark>analyze</mark> , markup language	analyze	>	w2	
W = AVG(w1, w2, w3, w1, w4, w2) If w1 is the nearest to W the chosen relation will be: semantic web, combine, markup language				

Triples Normalization and Merging

- The same relation may be expressed by different triples depending on the verb connecting the entities
- A predicates taxonomy has been built using word embeddings on lacksquarethe Microsoft Academic Graph and hierarchical clustering, for all possible relation labels

Every relation is associated with the corresponding word embedding of the verb expressing the relation

The relation with the word embeddings nearest to the average of all word embeddings is chosen as the most representative relation











- CSO Triples Integrator: it includes some triples derived with limited inference on CSO:
- given a triple (e1, r, e2), if in CSO the entity e3 is *superTopicOf* of e1, we also infer the triple (e3, r, e2).
- For instance, given the triple («NLP systems», «use», «Dbpedia»)
 - If hyp(«Semantic Web Technologies», «Dbpedia»), then we can infer the triple («NLP systems», «use», «Semantic Web Technologies»)

Triples Selection and Enrichment

- All triples coming from the Neural Extractor and OpenIE after
- the normalization and merging phases are kept
- Triples coming from CSO + « linking » verbs are kept only if
- their support is large enough (at least seen in 10 documents)







- - 87,030 from the Extractor Framework (TEF)
 - 8,060 from OpenIE (TOIE)
 - 14, 015 from the CSO + linking verbs method (TPoS).

Experiments

26,827 abstracts about the Semantic Web domain from Microsoft Academic Graph

The resulting knowledge graph contains 109,105 triples





semantic web	extends	WWW	
semantic web	supports	information integration	
semantic web	relies on	ontology	
ontology	provides	semantics	
ontology	represents	domain knowledge	
RDF	is	data model	
SPARQL	is	query language	
OWL	is	ontology language	
OWL	employs	open world assumption	
W3C	recommends	ontology	
Protégé	is	ontology editor	
Protégé	creates	OWL ontology	
semantic web services	extends	web services	

Example of extracted triples









Examples of extracted sub-graphs







- We built a Gold Standard composed of 818 triples
 - 401 triples from the Neural Extractor
 - 102 from Stanford Open IE
 - 170 triples from the CSO + linking verbs method
 - 212 randomly selected triples that were discarded by the framework pipeline.
- sense when compared with their knowledge on the subject) or not.
 - The agreement between experts was 0.747 ± 0.036

Evaluation

• Five experts in the field of Semantic Web annotated the triples as correct (that is, if they make







Approach	Precision	Recall	F-Score
Neural Extractor	0,72	0,55	0,62
OpenIE	0,65	0,13	0,21
CSO + linking verbs	0,73	0,24	0,36
Neural + OpenIE	0,70	0,66	0,68
Full Pipeline	0,70	0,81	0,75

Evaluation







•

•

•

Trivial knowledge (high-level entities with general relations) either: Does not appear (so trivial that it is not necessary to include into the paper) 0 Appears too often (outweigh all other relationships) In any case most (useful) relations appear only once Need to introduce a confidence weight over the knowledge when we extract it Is the sentence complex or simple? 0 Is the paragraph in the introduction/conclusions? 0 Some parameters may be specific to the considered domain Need more experiments!

Mining Scholarly Data for Scientific Knowledge Graph Construction

Conclusions







Acknowledgements

- Diego Reforgiato, Ass.Prof. Università di Cagliari
- Danilo Dessì, PhD Student Università di Cagliari
- Francesco Osborne, Research Fellow Knowledge Media Institute, The Open University, Milton Keynes



gliari liari edge Media Institute



