

Google Summer of Code, 2014

Progress Report

Organization: DBpedia

Project Name: Abbreviation Base – A multilingual knowledge base for abbreviations

Mentor Name: Martin Brümmer

<u>Week</u>	<u>Implementation</u>	<u>Progress from previous week</u>	<u>Problems</u>
Week 1 (April 28 - May 4)	<p>Implemented sample lemon model for 4 abbreviations (from abbrevs_en.tsv in the github repository) in RDF/XML format with the following specification:</p> <ul style="list-style-type: none"> • Every language has its own lemon:Lexicon • In every Lexicon, there are many lemon:LexicalEntry resources for the abbreviations. • Every LexicalEntry has a lemon:LexicalForm. This will contain the abbreviation string. • The LexicalEntry also has many lemon:LexicalSense resources. These contain the meaning of the abbreviations as a string and the link to the DBpedia page of the abbreviation as reference. • The LexicalEntry also has a translation link that links the LexicalSense resource to its equivalent in English, if it exists. <p>The above model was validated using 'rapper-utils' and was found syntactically correct.</p> <p>The correctness of a lemon file was tested with lemon validator: https://github.com/jmccrae/lemon-model.net/tree/master/validator and was found completely valid.</p>		Had some namespace problems which were resolved
	Similar implementation of abbreviation (4 abbreviations) was also done in turtle format with the above specification.		The correctness of a lemon file was tested with lemon

	<p>The above model was validated using 'rapper-utils' and was found syntactically correct.</p> <p>[the .ttl file has been pushed to the github repository: https://github.com/der-bruemmer/abbreviation-base]</p>		<p>validator: https://github.com/jmccrae/lemon-model.net/tree/master/validator and following error was found:</p> <p>[ERROR] Failed to parse as RDF: file:///home/div/Desktop/abbreviation_lemon_model.ttl:1:0: not well-formed (invalid token)</p>
<p>Week 2 (May 5 - May 11)</p>	<p>Made a python script.</p> <p>It did following things:</p> <ol style="list-style-type: none"> 1. Downloaded the required files from DBpedia dump data for a language ('en' was tested), read them and saved them in .bz2 format. 2. Decompressed the zipped files saved in hard disk and then according to the regular expression, writes a .ttl file containing the triples ending with ., ? or ! (as required). I used the bz2 function BZ2File(file_name,'rb'), to decompress the files. 	<p>Made a python script which reads necessary files from DBpedia server and decompresses them and store them in hard disk.</p>	<p>The file abbrev_en.ttl in the repo had many triples that were not present in http://downloads.dbpedia.org/3.9/en/redirects_en.ttl.bz2 . When comparing the file I have written via code, I found that triples after "<http://dbpedia.org/resource/American_Telephone_and_Telegraph_Inc> <http://dbpedia.org/ontology/wikiPageRedirects> <http://dbpedia.org/resource/AT&T_Corporation> ." are not present the redirects file that I had decompressed.</p>
<p>Week 3 (May 12 - May 19)</p>	<p>Improved the above code to automate the downloading and decompressing of files.</p> <p>I executed the code for French (fr), English (en) and also made a time file with the total time taken by each file to download and decompress.</p>	<p>I used os.popen library of python to decompress the file. The problem of incomplete file reading was solved.</p>	

	<p>The script downloads 7 files namely:</p> <ol style="list-style-type: none"> 1. redirects- contains abbreviations 2. instance_types_heuristic 3. interlanguage_links 4. interlanguage_links_chapters 5. labels 6. article_categories 7. disambiguations <p>All the things were automated.</p>	<p>Now all the files were downloaded successfully and decompressed completely.</p> <p>The problem was due to a bug in bz2 python library due to which I was unable to read complete file.</p>	
<p><u>First Phase Completed</u></p>	<ul style="list-style-type: none"> • Made the code to automate the downloading and decompressing of necessary files from DBpedia data dump. • I executed the code for French (fr), English (en) and also made a time file with the total time taken by each file to download and decompress. • The script downloads 7 files namely: <ol style="list-style-type: none"> 1. redirects- contains abbreviations 2. instance_types_heuristic 3. interlanguage_links 4. interlanguage_links_chapters 5. labels 6. article_categories 7. disambiguations <p>All the things were automated.</p>		

<p>Week 4 (May 20 - May 26)</p>	<p>Importing the necessary files into local virtuoso.</p> <p>(Setting up files into virtuoso so that it could be accessed using SPARQL queries)</p> <p>Setting up isql console to do the same.</p>		<p>I always get following error on executing: isql -v 1111 dba dba "[IM002][unixODBC] [Driver Manager]Data source name not found, and no default driver specified [ISQL]ERROR: Could not SQLConnect"</p> <p>Whenever I upload a big file around 2 GB in quad store upload in Virtuoso, I am automatically logged off and then I have to login again. It is not happening when I am uploading a smaller file.</p>
---	--	--	--

<p>Week 5 (May 27 - June 2)</p>	<p>Importing the necessary files into local virtuoso .</p> <p>Upload the file via the isql-vt console and execute queries on it.</p>	<p>Virtuoso may not use the standard isql command but either isql-vt or isql-v or isql-t.</p> <p>So I tried isql-vt and the isql front-end opened.</p> <p>Files were uploaded successfully as I zipped big files which took a lot of time (as VOS can import files in gzip format).</p>	<p>The code took about 3 hours to run and import files.</p>
<p>Week 6 (June 3 - June 9)</p>	<p>Working with the extractor code, extending the code already present in repository, heavily changing the SPARQL queries to extract the required information and improving the code to store the extracted information for the abbreviation.</p>	<p>Files were imported successfully.</p> <p>Tested SPARQL queries over uploaded data, the query is working perfectly fine for the abbreviation data.</p>	<p>Type and sameAs has multiple values for one abbreviation, so how is it to be fetched and how is to be stored and how will they be stored in Lemon file ?</p> <p>Eg: select * where {<http://dbpedia.org/resource/BC_Card> rdf:type ?type } (sense of <http://dbpedia.org/resource/BC> abbreviation B.C.) has following rdf:type</p> <p>type http://www.w3.org/2002/07/owl#Thing . http://schema.org/Organization http://dbpedia.org/ontology/Agent http://dbpedia.org/ontology/Company http://dbpedia.org/ontology/Company</p>

			<p>dbpedia.org/ontology/Organisation</p> <p><u>Solution as per me:</u> I will use a list to store the multiple values (type or sameAs) and will append in the existing dictionary in the code which I have already made.</p> <p>Problem in fetching category and rdf:type of the abbreviations which do not have disambiguations.</p>
<p>Week 7 (June 9 - June 15)</p>	<p>Fetching the necessary details of abbreviations like label, type, sameAs, categories.</p> <p>Multiple values like owl:sameAs, rdf:type, are stored in a list and stored it in the final list i.e. list in list. The format is :</p> <p>->for disambiguations: abbrevs[abbrev+" "+str(count)] = [result["label"]["value"], result["o"]["value"], list_sameAs, list_type]</p> <p>->for non-disambiguations: data = [result["label"]["value"], uri, list_sameAs]</p> <p>Data Fetched: #---result["name"]["value"] ==> lemon:value "Michel van der Aa"@en #---result["label"]["value"] ==> rdfs:label "Michel van der Aa" #---result["o"]["value"] ==> lemon:reference <http://dbpedia.org/resource/Michel_van_der_Aa></p>	<p>Extended the code, made some queries to extract the necessary information related to the abbreviations like rdf:type, owl:sameAs, category and label.</p> <p>Used list to store multiple values and stored this list into a final list making a 2D list.</p> <p><u>Solution to multiple sameAs and type value by mentor:</u> We don't really need owl#Thing, but we want "dbpedia.org" or "schema.org") excludes any "yago" types, that are somewhat messy.</p> <p>I have used REGEX in the query which</p>	

	<p>#----list_sameAs ==> owl:sameAs #----list_type ==> rdf:type "but works only for abbreviations with disambiguations" #----abbrevString ==> writtenRep , stores abbrev</p> <p>The code is ready, it fetches the required information and stores it in a .tsv file. I have tested the code for first five inputs and it works perfectly fine for both abbreviation with disambiguation and without it. (categories and type will be fetched later.)</p>	<p>will only fetch the required data i.e. dbpedia and schema.</p> <p><u>Solution to fetch category:</u> The category file being uploaded was incorrect. The correct file is "article_categories.ttl"</p> <p>They are related using: http://purl.org/dc/terms/subject</p>	
<p>Week 7 (June 16 - June 22)</p>	<p>Extend the above code to make a lemon model.</p> <p>End of First Phase:</p> <ul style="list-style-type: none"> ● download and extract the necessary files ● made a .tsv file ● make a lemon file using the code 	<p>I have extended the code to make lemon file. The code now makes the lemon file (not from the .tsv file, it makes it independently without using the .tsv file).</p> <p>I have tested the lemon file with rapper-utils, there is just one syntax error.</p>	<p>I have tested the lemon file with rapper-utils and here is the error report: rapper -gc abbreviation_lemon_en.ttl rapper: Parsing URI file:///home/div/Desktop/lemon%20file/abbreviation_lemon_en.ttl with parser guess rapper: Guessed parser name 'turtle' rapper: Error - URI file:///home/div/Desktop/lemon%20file/abbreviation_lemon_en.ttl:296 - syntax error at ' _' rapper: Parsing returned 867 triples</p> <p>Actually the problem is due to links like : "<http://it.dbpedia.org/resource/A.I._Intelligenza_artificiale>," which contain ' _'</p>

			<p>and I can not remove this '_' as then the link will become invalid. When I removed such links there is no error:</p> <pre> rapper -gc abbreviation_lemon_en.ttl rapper: Parsing URI file:///home/div/Desktop/lemon%20file/abbreviation_lemon_en.ttl with parser guess rapper: Guessed parser name 'turtle' rapper: Parsing returned 1060 triples </pre>
<p>Week 8 (June 23 -June29)</p>	<p>Deploy the software on the server, that is:</p> <p>Set up a script that:</p> <ul style="list-style-type: none"> -automatically downloads language specific DBpedia files needed by the extraction for one language -automatically imports these files into virtuoso -runs the extraction on it and exports lemon and tsv files -clears the graph in virtuoso after the extraction has ended -deletes the dbpedia files downloaded in the beginning -loads the next language from a list of languages to convert and start again at 1 <p>Like this it will automate the extraction.</p>	<p>setting up the server(with ssh keys) and importing the files on the server</p>	
	<p>Some changes in the code extractor_lemonMaker.py (by mentor):</p> <ul style="list-style-type: none"> -Use SPARQL contains function instead of a regex as regex is a much more costly (and slow) operation. -There is some hard coded "langMatches("EN")" in queries. -remove the dependence 		

	<p>on "only_roman_chars"</p> <p>-change the Namespace of lemon to "http://lemon-model.net" instead of monnetproject.eu.</p> <p>-change "@prefix : http://www.example.org/lexicon" to "@prefix : http://nlp.dbpedia.org/abbrevbase"</p> <p>-Lexicons should have uris of the form http://nlp.dbpedia.org/abbrevbase/lexicon/\$lang</p>		
<p>Week 9 (June 30 - July 6)</p>		<p>I have done the changes in the code as required.</p> <p>-corrected langMatches</p> <p>-removed the dependence on "only_roman_chars".</p> <p>-changed the Namespace of lemon to "http://lemon-model.net" instead of monnetproject.eu.</p> <p>-changed "@prefix : http://www.example.org/lexicon" to "@prefix : http://nlp.dbpedia.org/abbrevbase"</p> <p>-changes the querying owl:sameAs and rdf:type from regex to contains.</p> <p>-changed the locations of input and output files</p>	

		and tested. also fetched the categories successfully.	
Week 10 (July 7 - July 13)	<p>Executed the code for 'en' in the server. The final files (i.e lemon and tsv) have been created successfully for 'en' and they are present in /home/akswadmin/abbrev_extracted/en . The code is working perfectly fine now and is also automated.</p> <p>This script generates a report file server_report in /home/akswadmin/abbrev_repo. This is a report for 'en' :</p> <p>____Server Running Report____</p> <p>*****Current Language: en*****</p> <p>-----Files for en downloaded successfully-----</p> <p>-----Importing files into virtuoso-----</p> <p>-----</p> <p>-----Files imported-----</p> <p>-----Running lemon maker-----</p> <p>-----Removing graph-----</p> <p>-----Done en-----</p> <p>=====</p> <p>=====</p> <p>=====</p> <p>en took 12408 seconds</p> <p>=====</p> <p>=====</p> <p>=====</p>	<p>Executed the code for 'en' in the server. which makes lemon and tsv files.</p>	<p>For other languages like 'fr' and 'hi', there is a problem,</p> <p>1. there is no file like instance_types_heuristic_fr. ttl.bz2 in fr (and same goes for other languages) and some even don't have instance_types like http://downloads.dbpedia.org/3.9/bat_smg/ .</p> <p>2. Some do not have interlanguage_links_chapters. ttl.bz2 like http://downloads.dbpedia.org/3.9/bat_smg/</p>
	<p>Other corrections in the code:</p> <p>-Space in @prefix rdf definition</p> <p>-Missing ";" after "a lemon:Lexicon"</p> <p>-&nbsp; should be changed to " _"</p> <p>-"lemon:entry" triples don't match the entry resources in the file:</p>		<p>Found that some abbreviations are invalid in lemon</p> <p>1. Those starting with digits (though they can contain digits)</p> <p>2. Those containing</p>

	<p>\$Lexicon lemon:entry :AD_Entry but <A.D._Entry> a lemon:LexicalEntry . should be</p> <p>\$Lexicon lemon:entry :A.D._entry :A.D._entry a lemon:LexicalEntry .</p>		<p>symbols like =, (), %, ', &, ... , & etc</p> <p>There is also a problem in the extraction code in regard to URL formatting. There is need to do is append a "<" and a ">" at the beginning and the end when added as objects to triples instead of string replace. It introduced problems in the RDF, because URLs can contain "," and by doing replace(",",">").</p> <p>http://dbpedia.org/resource/Category:University of California>,_Irvine_alumni></p> <p>Instead, it should be the single URLs and append the brackets before adding it to the string.</p>
	<p>- Updated the code to replace : lemon:entry :AD_Entry to lemon:entry :A.D._entry , :A.I.entry, ... i.e. in their original form</p> <p>-Made a function that will append "<", ">" as required</p>	<p>The strings should not be contained in any URI if format them like :A.D._entry. Besides that, digits, parenthesis and ampersands are allowed in URIs</p> <p>For the missing files, accept that there might be data missing. Most important are the abbreviation strings and their full meaning. So included exception</p>	<p>I. On replacing :AD_Entry with :A.D._Entry, the rapper gives a syntax error. else no error.</p> <p>II. If there is any symbol in :AD_Entry like :AD%_Entry, :A!D_Entry etc then also rapper gives a syntax error. There are several abbreviations which contain symbols as like :%3F_Entry , :%5C_Entry , :Negima%3F_Entry</p>

		error handling in the code.	<p>III. If there is any entry starting with a digit then also rapper gives a syntax error like :0_Entry</p> <p>IV. On changing <A.D._Entry> a lemon:LexicalEntry . to :A.D._entry a lemon:LexicalEntry . then also there is a syntax error.</p>
	<p>-Executed the latest code in the server for just 5 abbreviations and the result is as follows:</p> <pre> root@abbreviationBase:~/en# rapper -gc abbreviation_lemon_en.ttl rapper: Parsing URI file:///home/akswadmin/ en/abbreviation_lemon_en.ttl with parser guess rapper: Guessed parser name 'turtle' rapper: Parsing returned 2806 triples </pre>	<p>Outdated rapper, so files verified on server and the rapper now gives a correct report which clears the problems I, III and IV.</p> <p>Using full URIs</p> <pre> <http:// nlp.dbpedia.org/ abbrevbase/lexicon/ en> a lemon:Lexicon ; lemon:language "en" ; lemon:entry . <http:// nlp.dbpedia.org/ abbrevbase/lexicon/ en/entry/A.D.> . <http:// nlp.dbpedia.org/ abbrevbase/lexicon/ en/entry/A.D.> lemon:form . <http:// nlp.dbpedia.org/ abbrevbase/lexicon/ </pre>	<p>Tested on a similar file, added %, parenthesis and ampersand in entry elements (filename: /home/akswadmin/en/abbr_lemon_experiment.ttl) and here is the rapper result:</p> <pre> root@abbreviationBase:~/en# rapper -gc abbr_lemon_experiment.ttl rapper: Parsing URI file:///home/ akswadmin/en/ abbr_lemon_experiment.ttl with parser guess rapper: Error - URI file:///home/ akswadmin/en/ abbr_lemon_experiment.ttl:12 - syntax error at '%' rapper: Guessed parser name 'turtle' rapper: Failed to parse file abbr_lemon_experiment.ttl guess content rapper: Parsing returned 3 triples </pre>

	<p>ran the code for 1000 abbreviations and the lemon file I received was completely correct.</p> <p>The code has been executed for en and the lemon file is perfect as per the rapper.</p> <p>here is the rapper validation:</p> <pre> root@abbreviationBase:~# rapper -gc en/ abbreviation_lemon_en.ttl rapper: Parsing URI file:///home/akswadmin/ en/abbreviation_lemon_en.ttl with parser guess rapper: Guessed parser name 'turtle' rapper: Parsing returned 1663426 triples </pre>	<pre> en/entry/A.D.#form> ; lemon:sense <http:// nlp.dbpedia.org/ abbrevbase/lexicon/ en/entry/ A.D.#sense> ; a lemon:LexicalEntry . <http:// nlp.dbpedia.org/ abbrevbase/lexicon/ en/entry/A.D.#form> lemon:written Rep "A.D."@en ; a lemon:LexicalForm . <http:// nlp.dbpedia.org/ abbrevbase/lexicon/ en/entry/A.D.#sense> lemon:definiti on [lemon:value " Anno Domini"@en] ; lemon:referen ce <http:// dbpedia.org/ resource/ Anno_Domini> ; a lemon:LexicalSense . This solves the issue as chars and digits are allowed in URIs. </pre>	<p>(also checked individually for %, (), & and , and all of them resulted into syntax error.)</p> <p>-An error due to lemon:value "Oran " Juice" Jones"@en . Its due to multiple " ", so corrected it using "Oran \"Juice\" Jones"@en</p> <p>-On executing the code for all the abbreviations then there might be some error due to which the code stops</p> <p>-Problem solved using exception handling:</p> <p>Index Error when converting type, category etc to string. The error was obtained when any of them were empty i.e. when the query did not return the result for them, so I used exception handling to remove it.</p> <p>Then there were errors due to not initializing variables <code>rdf_type</code>,</p>
--	---	---	---

			label, category, owl_sameAs, ref which gave an error when again the query was not able to return the results to the respective query.
<u>Second Phase Completed</u>	As a next task, run the script for de, nl, fr, es, it, ru, cs, pl, sv, da, et, pt, no, nn, sl, fi		
Week 11 (July 14 - July 20)		completed extraction for 'en', the lemon file made is absolutely correct as per rapper validation. Executed the code for other languages in the database.	getting significant differences in the count of abbreviations between my extraction and the extraction mentor originally did. For German, tsv has 4689 lines, mine only has 696. For French and Dutch, you extract around 300 lines less.
	Get the URIs right Abbreviations containing "_" should not be included at all In the Polish and French abbreviations, there are URL encoding (%3F instead of ?). If encounter "%" in the abbreviations, run a URLdecode function over it.	The problem of incomplete abbreviations was that for those abbreviations whose query returned nothing, the dictionary abbrevs was no formed, thus it was not written in the file. I just included a condition which prevents the above problem. http://localhost:8890/sparql	there are abbreviations in there that contain "_", thus massively increasing the scope. For example, in the German file, there are a lot of dates. Also something does not work right with the URIs: 2980er_v._Chr. 30. Jahrhundert v. Chr. "30. Jahrhundert v. Chr."@de <ttp:// dbpedia.org/resource/30. Jahrhundert v. Chr >

		<p>query this endpoint for every language, to remove problem of missing data.</p>	<p>missing "h" at the start of the URI and the missing period at the end of the URI. From what I have seen, it only happens to URIs that contain "_".</p> <p>missing some disambiguations for the abbreviation "A.A." in the German file. Disambiguations.</p>
	<p>I produced an overall 22859 abbreviations with 78197 meanings in 99 languages. Overall unique abbreviations (aggregated over the languages) are 15646.</p> <p>The other languages (20 out of 119) for which abbreviations could not be extracted had either inappropriate data, or data not desired.</p> <p>deleted the old extractions, the dbpedia files etc</p> <p>clear the virtuoso graph, make a new graph , "http://nlp.dbpedia.org/abbrevbase/" and import all abbreviations from all languages into it.</p>	<p>Removed all the abbreviations which only have digits and no character.</p>	<p>we have downloaded /instance_types_heuristic_\$language.ttl.bz2 for instance types. But some languages do not have instance_types_heuristic_\$language.ttl.bz2 rather they have instance_types_\$language.ttl.bz2 so we need to download both the files (which ever available) for all the languages.</p>
<u>Third Phase Completed</u>	<p>I produced an overall 22859 abbreviations with 78197 meanings in 99 languages. Overall unique abbreviations (aggregated over the languages) are 15646.</p>		