

INTERNATIONAL INSTITUTE OF INFORMATION TECHNOLOGY, HYDERABAD

CS4.501 SOCIAL COMPUTING SPRING-2021

COURSE PROJECT - SECOND DELIVERABLE

Project Outline (final)

**Enhance Knowledge Graphs/Knowledge bases
like Wikidata**

Team Members:

Pushpa Yadhav (2019900034)
Vijaya Lakshmi (2018900071)
Karthikeyan Arumugam(2018900074)

Mentor:

Tushar Abhishek

April 1,2021



INTERNATIONAL INSTITUTE OF
INFORMATION TECHNOLOGY

H Y D E R A B A D

1 Update Since Last Submission

- Quality of generated triplets from Match algorithm was not great in our last submission, Majorly due to coreference resolution challenges. Eg often in large documents, Subject reference which happens after first occurrence will have third person/second person notation, hence generated Triplets often had Subject/Object entity as He , She, they etc.
- to over come this issue - we had to add additional preprocessing step to the existing clean up process to add coreference resolution to it.
- Coreference resolution issues in Pre Processing - 85 % completed
 - Hugging face - Core ref Approach - Model yields average performance, in many cases it is not resolving coreferences effectively and it is getting struck with one Object/Subject name and retains it for rest of the document

eg

	subject	relation	object	st
0	American actress	released	American actress	N
1	American actress	include	The Singing Nun	N
2	Reynolds	reached	new younger generation	PI
3	American actress	released	American actress	N
4	American actress	released	second autobiography	N
5	Reynolds	had	business ventures	PI
6	Reynolds	received	the Screen Actors Guild Life Achievement Award	PI
7	American actress	received	the Screen Actors Guild Life Achievement Award	N
8	Debbie	has	refreshing sense	PI
9	Debbie	led to	role	PI
10	Reynolds	reversed	Reynolds opinion	PI
11	Debbie	portrayed	Jeanine Deckers	PI
12	Debbie	played	Helen Chappel Hackett's mother	PI
13	Debbie	played	recurring role	PI
14	Debbie	had	cameo role	PI
15	Reynolds	scored	hits	PI
16	Debbie	released	album	PI
17	production	broke	records	N
18	Debbie	received	Tony nomination	PI
19	replacement	own	West End	N
20	2008	displayed	items	D
21	Reynolds	opened	Reynolds own dance studio	PI
22	Debbie's Wayf. Reynolds	purchased	the Clarion Hotel and Casino	PI
23	The Eddie Fisher – Elizabeth Taylor affair	noted	bright	O

Hence we didnt used hugging face approach any further.

- AllenNLP - Coreference resolution yields better results than hugging face neural coref library, Challenge with size of the ip document , we ran into issues when input document is greater than 1000 characters long.
- Map - Algorithm updates
 - In Last submission we were able to provide entity linking, hence generate triplets were mapped to existing KG entities. but we had challenge with Resolving relationship/edges text to relation ID in wiki data

- We employed SPARQL query to collect list of preexisting relations in the Wiki data. used word similarity to identify if generated relationship/edge is preexisting entry in wikidata
- Relationship resolution / Edge resolution to existing entries in Wiki-data : Partly completed.- 50% , Few more tweaking and changes are required to make it more functional
- Confidence Score for the generated Triplets are computed using $P(\text{generated triplets})$ in language model. More fine tuning is required to build a confidence score measure.
-

2 Work In progress Action Items

- Running the entire pipeline for large volume of data - 40% completed
- Qualitative evaluation - 20 %
- Metrics to measure increase in density of the KG - 20%
- relationship resolution using SPARQL query - Fine tuning
- confidence score of generated triplets - Fine tuning
- Qualitative evaluation for manually comparing a single wiki page - to gauge accuracy of mining KG graph entries from the proposed method.

3 Remaining Action Items and Target dates

- Running entire data pipeline for Already scrapped large corpus - 4th April
- Qualitative Evaluation of the Sample generated KG - 6th April 2021
- Quantitative metric results - 8th April 2021
- Final Presentation and reports - 10th April 2021

4 References

Git Hub Link: <https://github.com/Team9-CS4-501-S21-CourseProject/EnhanceWikiData/>