

Attribute Template Generation

Daksh, Vasu, Divanshi

Task

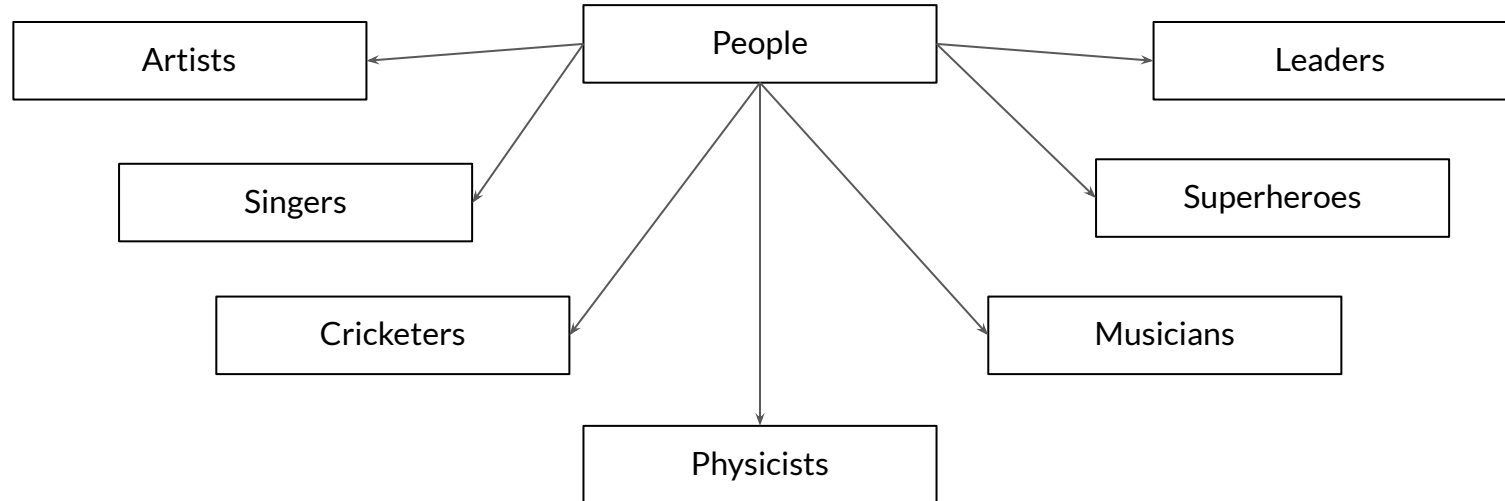
- Given a specific domain, and a set of input documents, the task is to generate an attribute template corresponding to the input domain which can be used in generating new articles for that domain.
- Secondary task - Once the list of attributes for the domain is extracted, we need to divide them in accordance with the section of the wikipedia page a certain group of attributes would belong to.

Scope

- Input
 - Domain name
 - ~1K wikipedia articles
 - Language - English
- Attributes to be extracted from
 - Infobox
 - Wikidata
 - Wikidata page
 - Wikipedia text
- Output
 - Template for that domain
 - Intermediate - List of attributes and their categorization for the domain

Domain

- For given wikipedia domain, a large number of sub categories exist which may cover very different among themselves. So, instead of focusing on a domain, it will be better to start off by focusing on a sub-domain.



Sub-category stats

The below table shows stats of some sub-categories we manually picked. Out of these, we have decided to work on the 4 highlighted categories for now. We have mainly made the judgement based on number of distinct infobox properties as well as the average size of an infobox for each domain

Category	Depth	No of articles	Sum of infobox sizes	Average infobox size	Number of distinct Properties
Music directors	2	2492	9315	3.737961477	229
Artists	1	1792	8233	4.594308036	570
Mathematicians	1	2161	12819	5.931975937	324
Physicists	1	1073	6738	6.279589935	329
Novelists	1	1736	11598	6.680875576	301
Writers	1	2372	15993	6.742411467	587
Actors	1	1091	7864	7.208065995	512
Musicians	1	2467	20347	8.247669234	543
Feminists	1	2538	22334	8.799842396	579
Singers	1	823	7616	9.253948967	386
Supervillains	2	1969	19121	9.711020823	326
Superheroes	2	2528	31475	12.4505538	593
Leaders	1	1440	23817	16.53958333	574
Politicians convicted of crimes	2	1270	21463	16.9	546
Cricketers	1	926	26345	28.45032397	381

Steps

- First step would be to simply extract the already present attributes in the infobox and wikidata for a given wikipedia page.
- To get attributes from the text
 - Use an OpenIE system to extract relations from the text.
 - We will look at the shortcoming of existing openIE systems and try to come up with heuristics to resolve them ourselves.