

KG/pre-training resource/corpus

数据集描述&例子&统计信息	使用数据集名称	数据集出处	使用数据集论文	数据集链接	数据集处理代码	备注																																
<p>Wikipedia abstract数据，abstract里面有人工标注的实体</p> <hr/> <p><i>The Lord of the Rings</i> is an <i>epic high-fantasy</i> novel written by English author <i>J. R. R. Tolkien</i>.</p> <hr/> <p>Entity Annotations: The Lord of the Rings, Epic (genre), High fantasy, J. R. R. Tolkien</p> <hr/> <p>4.4million的Wikipedia articles</p> <table><tr><th>Language</th><th>Abstracts</th><th>Entity links</th><th>Triples</th></tr><tr><td>Dutch</td><td>1,740,494</td><td>11,344,612</td><td>114,284,973</td></tr><tr><td>English</td><td>4,415,993</td><td>39,650,948</td><td>387,953,239</td></tr><tr><td>French</td><td>1,476,876</td><td>11,763,080</td><td>116,205,859</td></tr><tr><td>German</td><td>1,556,343</td><td>15,859,142</td><td>153,626,686</td></tr><tr><td>Italian</td><td>907,329</td><td>7,705,247</td><td>75,698,533</td></tr><tr><td>Spanish</td><td>1,038,639</td><td>11,558,121</td><td>111,293,569</td></tr><tr><td>All</td><td>11,135,674</td><td>97,881,150</td><td>959,062,859</td></tr></table> <p>包括信息：abstract，abstract中的surface form（mention），mention的position，mention对应的KG中的entity（一个链接），mention在语料中出现的次数</p> <p>{Berlin;http://dbpedia.org/resource/Berlin; 9338} indicates 9.338occurrences of the entity "Berlin" and includes the link tothe corresponding DBpedia resource.</p> <p>数据集包括两大文件：一个abstract内容，如第一个图；一个surface forms，如这个eg。</p> <p>它提供abstract和entity的对齐了没？</p> <p>Surface form是什么？ DBpedia和Wikipedia的 entity id能对上嘛？</p> <p>Entity 的size是多少？ 比abstract多吧？ surface form语料是2.6M。</p>	Language	Abstracts	Entity links	Triples	Dutch	1,740,494	11,344,612	114,284,973	English	4,415,993	39,650,948	387,953,239	French	1,476,876	11,763,080	116,205,859	German	1,556,343	15,859,142	153,626,686	Italian	907,329	7,705,247	75,698,533	Spanish	1,038,639	11,558,121	111,293,569	All	11,135,674	97,881,150	959,062,859	DBpedia abstract corpus	DBpedia Abstracts: A Large-Scale, Open, Multilingual NLP Training Corpus	Learning Distributed Representations of Texts and Entities from Knowledge Base, LUKE 大概也是用了这个语料，但是预训练的 entity vocab size是1M，有点少啊，而且论文对这一块介绍太少了。	http://download.dbpedia.org/2015-04/ext/nlp/abstracts/	http://github.com/studio-ousia/ntee	这个工作提到的 entity 应该是我们常说的 mention，也就是数据集论文中的 surface form，也就图片中划下划线的 mention
Language	Abstracts	Entity links	Triples																																			
Dutch	1,740,494	11,344,612	114,284,973																																			
English	4,415,993	39,650,948	387,953,239																																			
French	1,476,876	11,763,080	116,205,859																																			
German	1,556,343	15,859,142	153,626,686																																			
Italian	907,329	7,705,247	75,698,533																																			
Spanish	1,038,639	11,558,121	111,293,569																																			
All	11,135,674	97,881,150	959,062,859																																			

english数据。entity和description对齐数据,使用 Wikidata dump 2019和Wikipedia数据。description使用的是Wikipedia page中的first section。relation是从Wikidata中抽取的。具体，对于wikidata中的每一个entry，找到他的Wikipedia page，抽取first section作为description。如果找不到Wikipedia page或者page长度小于5个词，则discard。从wikidata中抽取relation时，只保留头尾实体都没有被discard的，并且该relation在wikidata中也有page的。

Dataset	#entity	#relation	#training	#validation	#test
FB15K	14, 951	1, 345	483, 142	50, 000	59, 071
WN18	40, 943	18	141, 442	5, 000	5, 000
FB15K-237	14, 541	237	272, 115	17, 535	20, 466
WN18RR	40, 943	11	86, 835	3, 034	3, 134
Wikidata5M	4, 594, 485	822	20, 614, 279	5, 163	5, 133

wikidata的QID或者entity name和Wikipedia的page是怎么对齐的?

wikidata5M
KEPLER
KEPLER
<https://github.com/tunlp/ER-NIE>
<https://github.com/tunlp/ER-NIE>

知识图谱KB和free text对齐数据
Wikipedia abstract和wikidata triple对齐数据，指的啥？？

Dataset	Documents / Format	Unique predicates	Aligned Triples	Available
NYT-FB	1.8M sent.	258	39K	partially
TAC KBP	90K sent.	41	122K	closed
Google-RE	60K sent.	5	60K	publicly
FB15K-237	2.7 M patterns	237	2.7M	publicly
Wikireadings	4.7M articles	884	n.a.	publicly

Table 1: Statistics on predicate alignments from random walk.

为抽取T-REx，设计了一套pipeline，document信息使用的DBpedia Abstract Corpus（就是上面提到的那个），triple信息使用的wikidata truthy dump（144M triple）

Annotator	Documents covered	Alignments	Numerical Alignments	Uniq predicates
NYT-FB	1.8M	39K	None	258
TAC-KBP	0.09M	122K	n.a.	41
T-REx-SPO	0.79M	1.2M	21K	336
T-REx-NoSub	2.85M	5.2M	561K	642
T-REx-AllEnt	3.09M	11.1M	350K	633

Table 2: Number of alignments in different datasets.

包括json和NIF两种格式

Json example

NIF example

T-REx
T-REx: A Large Scale Alignment of Natural Language with Knowledge Base Triples
K-Adaptor
<https://old.datashub.io/dataset/tunlp-er-nie>
或者<https://hadyelsahar.github.io/tunlp-er-nie/>
做了一些处理后再使用，产生了T-REx-rc（只包括relation的表面形式出现的sentence，丢掉少于50个entity对的relation，最

			终得到 430个 relatio n和 5.5M个 senten ce)			
--	--	--	--	--	--	--