

[Get unlimited access](#)[Open in app](#)

Published in Towards Data Science

You have **2** free member-only stories left this month. [Upgrade for unlimited access.](#)



Chris Thornton

[Follow](#)Feb 9, 2020 · 4 min read ★ · [Listen](#)

Save



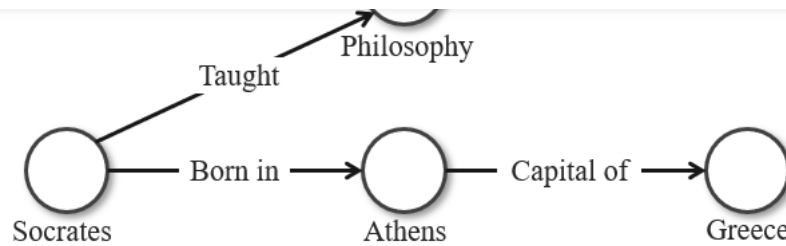
Auto-Generated Knowledge Graphs

Utilize an ensemble of web scraping bots, computational linguistics, natural language processing algorithms and graph theory.



Knowledge graphs are a tool of data science that deal with interconnected **entities** (people, organizations, places, events, etc.). Entities are the **nodes** which are connected via **edges**. Knowledge graphs consist of these **entity pairs** that can be traversed to uncover meaningful connections in unstructured data.





There are issues inherent with graph databases, one being the manual effort required to construct them. In this article I will discuss my research and implementations of automatic generation using web scraping bots, computational linguistics, natural language processing (NLP) algorithms and graph theory (with python code provided).

Web Scraping

The first step in constructing a knowledge graph is to gather your sources. One document may be enough for some purposes, but if you want to go deeper and crawl the web for more information there are multiple ways to achieve this using web scraping. Wikipedia is a decent starting point, as the site functions as a user-generated content database with citations to mostly reliable secondary sources, which vet data from primary sources.

Side Note: Always check your sources. Believe it or not, not all information on the internet is true! For a heuristic based solution, cross-reference other sites or opt for SEO metrics as a proxy for trust-signals.

I will avoid screen-scraping wherever possible by using a direct python wrapper for the [Wikipedia API](#).

The following function searches Wikipedia for a given topic and extracts information from the target page and its internal links.

```
1 import wikipediaapi # pip install wikipedia-api
2 import pandas as pd
3 import concurrent.futures
4 from tqdm import tqdm
5
6 def wiki_scrape(topic_name, verbose=True):
7     def wiki_link(link):
8         try:
9             page = wiki_api.page(link)
10            if page.exists():
11                return {'page': link, 'text': page.text, 'link': page.fullurl,
12                        'categories': list(page.categories.keys())}
13        except:
14            return None
15
16    wiki_api = wikipediaapi.Wikipedia(language='en',
17                                     extract_format=wikipediaapi.ExtractFormat.WIKI)
18    page_name = wiki_api.page(topic_name)
19    if not page_name.exists():
20        print('Page {} does not exist.'.format(topic_name))
21        return
22
23    page_links = list(page_name.links.keys())
24    progress = tqdm(desc='Links Scraped', unit='', total=len(page_links)) if verbose else None
25    sources = [{'page': topic_name, 'text': page_name.text, 'link': page_name.fullurl,
26               'categories': list(page_name.categories.keys())}]
27
28    with concurrent.futures.ThreadPoolExecutor(max_workers=5) as executor:
29        future_link = {executor.submit(wiki_link, link): link for link in page_links}
30        for future in concurrent.futures.as_completed(future_link):
31            data = future.result()
32            sources.append(data) if data else None
33            progress.update(1) if verbose else None
34    progress.close() if verbose else None
35
36    namespaces = ('Wikipedia', 'Special', 'Talk', 'LyricWiki', 'File', 'MediaWiki',
```

[Get unlimited access](#)[Open in app](#)

```
42     sources['topic'] = topic_name
43     print('Wikipedia pages scraped:', len(sources))
44
45     return sources
```

wikipedia_scrape.py hosted with ❤ by GitHub

[view raw](#)

Let's test this function on the topic: "Financial crisis of 2007–08"

```
wiki_data = wiki_scrape('Financial crisis of 2007-08')
```

Output:

Wikipedia pages scraped: 798

topic	page	text	link	categories
Financial crisis of 2007-08	Financial crisis of 2007-08	The financial crisis of 2007-08, also known as the global financial crisis a...	https://en.wikipedia.org/wiki/Financial_crisis_of_2007%E2%80%9308	['2000s economic history', '2007 in economics', '2008 in economics', 'All Wi...
Financial crisis of 2007-08	10-Q	Form 10-Q, (also known as a 10-Q or 10Q) is a quarterly report mandated by...	https://en.wikipedia.org/wiki/Form_10-Q	['Articles with short description', 'SEC filings']
Financial crisis of 2007-08	1970s energy crisis	The 1970s energy crisis occurred when the Western world, particularly the Un...	https://en.wikipedia.org/wiki/1970s_energy_crisis	['1970s economic history', '1979 in economics', 'All NPOV disputes', 'All ar...
Financial crisis of 2007-08	1973 oil crisis	The 1973 oil crisis began in October 1973 when the members of the Organizat...	https://en.wikipedia.org/wiki/1973_oil_crisis	['1973 in economics', '1973 in international relations', 'All articles ...
Financial crisis of 2007-08	1973-1975 recession	The 1973-1975 recession or 1970s recession was a period of economic sta...	https://en.wikipedia.org/wiki/1973%E2%80%931975_recession	['1970s economic history', '1973 in economics', '1974 in economics', '1975 i...
Financial crisis of 2007-08	1973-74 stock market crash	The 1973-74 stock market crash caused a bear market between January 1973 and D...	https://en.wikipedia.org/wiki/1973%E2%80%9374_stock_market_crash	['1973 in economics', '1974 in economics', 'Stock market crashes', 'Use...
Financial crisis of 2007-08	1973-75 recession	The 1973-1975 recession or 1970s recession was a period of economic sta...	https://en.wikipedia.org/wiki/1973%E2%80%931975_recession	['1970s economic history', '1973 in economics', '1974 in economics', '1975 i...
Financial crisis of 2007-08	1976 IMF crisis	The 1976 IMF Crisis was a financial crisis in the United Kingdom in 1976 w...	https://en.wikipedia.org/wiki/1976_IMF_crisis	['1976 in the United Kingdom', 'Economic history of the United Kingdom', 'Financi...
Financial crisis of 2007-08	1979 oil crisis	The 1979 (or second) oil crisis or oil shock occurred in the world due to dec...	https://en.wikipedia.org/wiki/1979_oil_crisis	['1979 in economics', '1979 in international relations', 'All articles ...

If you want to extract a single page use the below function:

```
1 def wiki_page(page_name):
2     wiki_api = wikipediaapi.Wikipedia(language='en',
3     extract_format=wikipediaapi.ExtractFormat.WIKI)
4     page_name = wiki_api.page(page_name)
5     if not page_name.exists():
6         print('Page {} does not exist.'.format(page_name))
7         return
8
9     page_data = pd.DataFrame({
10         'page': page_name,
```





```
16
17     return page_data
```

wikipedia_scrape_page.py hosted with ❤ by GitHub

[view raw](#)

Computational Linguistics & NLP Algorithms

Knowledge graphs can be constructed automatically from text using **part-of-speech** and **dependency parsing**. The extraction of entity pairs from grammatical patterns is fast and scalable to large amounts of text using NLP library **SpaCy**.

The following function defines entity pairs as **entities/noun chunks with subject — object dependencies connected by a root verb**. Other rules-of-thumb can be used to produce different types of connections. This kind of connection can be referred to as a **subject-predicate-object triple**.

```
1  import pandas as pd
2  import re
3  import spacy
4  import neuralcoref
5
6  nlp = spacy.load('en_core_web_lg')
7  neuralcoref.add_to_pipe(nlp)
8
9
10 def get_entity_pairs(text, coref=True):
11     # preprocess text
12     text = re.sub(r'\n+', '.', text) # replace multiple newlines with period
13     text = re.sub(r'[\d+]', ' ', text) # remove reference numbers
14     text = nlp(text)
15     if coref:
16         text = nlp(text._.coref_resolved) # resolve coreference clusters
17
18     def refine_ent(ent, sent):
19         unwanted_tokens = (
20             'PRON', # pronouns
21             'PART', # particle
22             'DET', # determiner
23             'SCONJ', # subordinating conjunction
24             'PUNCT', # punctuation
25             'SYM', # symbol
26             'X', # other
27         )
28         ent_type = ent.ent_type_ # get entity type
29         if ent_type == '':
30             ent_type = 'NOUN_CHUNK'
31             ent = ' '.join(str(t.text) for t in
32                             nlp(str(ent)) if t.pos_
33                             not in unwanted_tokens and t.is_stop == False)
34         elif ent_type in ('NOMINAL', 'CARDINAL', 'ORDINAL') and str(ent).find(' ') == -1:
35             refined = ''
36             for i in range(len(sent) - ent.i):
37                 if ent.nbor(i).pos_ not in ('VERB', 'PUNCT'):
38                     refined += ' ' + str(ent.nbor(i))
39             else:
```



```
45 sentences = [sent.string.strip() for sent in text.sents] # split text into sentences
46 ent_pairs = []
47 for sent in sentences:
48     sent = nlp(sent)
49     spans = list(sent.ents) + list(sent.noun_chunks) # collect nodes
50     spans = spacy.util.filter_spans(spans)
51     with sent.retokenize() as retokenizer:
52         [retokenizer.merge(span, attrs={'tag': span.root.tag,
53                                         'dep': span.root.dep}) for span in spans]
54     deps = [token.dep_ for token in sent]
55
56     # limit our example to simple sentences with one subject and object
57     if (deps.count('obj') + deps.count('dobj')) != 1\
58         or (deps.count('subj') + deps.count('nsubj')) != 1:
59         continue
60
61     for token in sent:
62         if token.dep_ not in ('obj', 'dobj'): # identify object nodes
63             continue
64         subject = [w for w in token.head.lefts if w.dep_
65                   in ('subj', 'nsubj')] # identify subject nodes
66         if subject:
67             subject = subject[0]
68             # identify relationship by root dependency
69             relation = [w for w in token.ancestors if w.dep_ == 'ROOT']
70             if relation:
71                 relation = relation[0]
72                 # add adposition or particle to relationship
73                 if relation.nbor(1).pos_ in ('ADP', 'PART'):
74                     relation = ' '.join((str(relation), str(relation.nbor(1))))
75             else:
76                 relation = 'unknown'
77
78             subject, subject_type = refine_ent(subject, sent)
79             token, object_type = refine_ent(token, sent)
80
81             ent_pairs.append([str(subject), str(relation), str(token),
82                               str(subject_type), str(object_type)])
83
84     ent_pairs = [sublist for sublist in ent_pairs
85                  if not any(str(ent) == '' for ent in sublist)]
86     pairs = pd.DataFrame(ent_pairs, columns=['subject', 'relation', 'object',
87                                             'subject_type', 'object_type'])
88     print('Entity pairs extracted:', str(len(ent_pairs)))
89
90     return pairs
```

get_entity_pairs.py hosted with ❤ by GitHub

[view raw](#)



Call the function on the main topic page:

```
pairs = get_entity_pairs(wiki_data.loc[0, 'text'])
```

Output:

Entity pairs extracted: 71

subject	relation	object	subject_type	object_type
Dow Jones Industrial Average	hit	Dow Jones Industrial Average peak closing price	NOUN_CHUNK	NOUN_CHUNK
Bank of America	purchased	Merrill Lynch	ORG	ORG
The Federal Reserve	took over	American International Group	ORG	ORG
The Reserve Primary Fund	broke	buck	ORG	NOUN_CHUNK
Congress	passed	the Emergency Economic Stabilization Act	ORG	LAW
Two of the three Big Three automobile manufacturers	received	bailout	CARDINAL	NOUN_CHUNK
Congress	approved	the American Recovery and Reinvestment Act	ORG	ORG
Dow Jones	hit	Dow Jones lowest level	NOUN_CHUNK	NOUN_CHUNK
Low interest rates	encouraged	mortgage lending	NOUN_CHUNK	NOUN_CHUNK
implicit guarantee	created	moral hazard	NOUN_CHUNK	NOUN_CHUNK

Coreference resolution significantly improves entity pair extraction by normalizing the text, removing redundancies, and assigning entities to pronouns (*see my article on coreference resolution below*).

Coreference Resolution in Python

Integrate Neural Network-Based Coreference Resolution into your NLP Pipeline using NeuralCoref

towardsdatascience.com

It may also be worthwhile to train a [custom entity recognizer model](#) if your use-case is domain-specific (healthcare, legal, scientific).

Graph Theory

Next, let's draw the network using the **NetworkX** library. I will create a **directed multigraph** network with nodes sized in proportion to **degree centrality**.

```
1 import networkx as nx
```



[view raw](#)

[Get unlimited access](#)[Open in app](#)

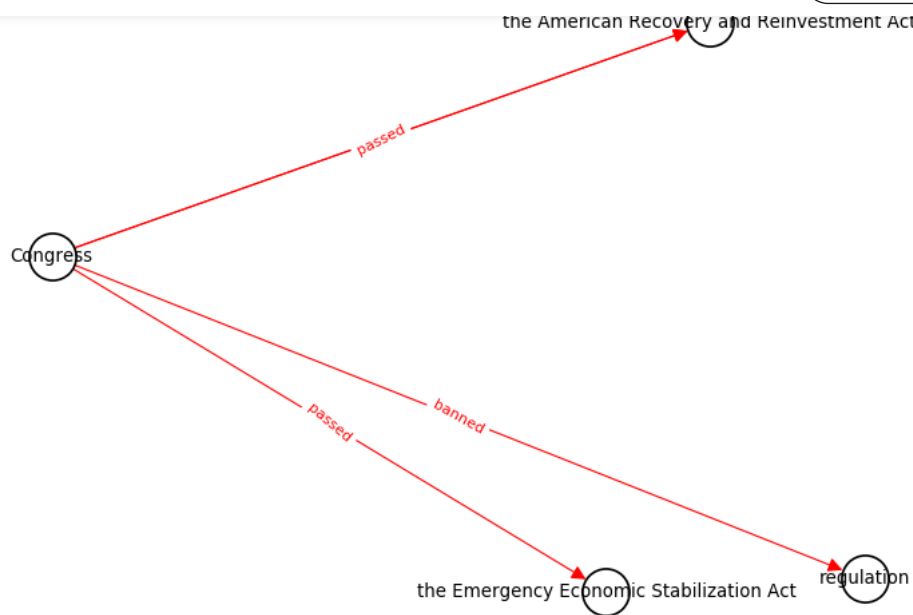
```
3     create_using=nx.MultiDiGraph))
4     edges = nx.dfs_successors(k_graph, node)
5     nodes = []
6     for k, v in edges.items():
7         nodes.extend([k])
8         nodes.extend(v)
9     subgraph = k_graph.subgraph(nodes)
10    layout = (nx.random_layout(k_graph))
11    nx.draw_networkx(
12        subgraph,
13        node_size=1000,
14        arrowsize=20,
15        linewidths=1.5,
16        pos=layout,
17        edge_color='red',
18        edgecolors='black',
19        node_color='white'
20    )
21    labels = dict(zip((list(zip(pairs.subject, pairs.object))),
22                      pairs['relation'].tolist()))
23    edges= tuple(subgraph.out_edges(data=False))
24    sublabels ={k: labels[k] for k in edges}
25    nx.draw_networkx_edge_labels(subgraph, pos=layout, edge_labels=sublabels,
26                                font_color='red')
27    plt.axis('off')
28    plt.show()
```

filter_graph.py hosted with ❤ by GitHub

[view raw](#)

```
filter_graph(pairs, 'Congress')
```



[Get unlimited access](#)[Open in app](#)

Knowledge Graphs at Scale

To effectively use the entire corpus of ~800 Wikipedia pages for our topic, use the columns created in the `wiki_scrape` function to add properties to each node, then you can track which pages and categories each node lies in.

I recommend using **multiprocessing** or **parallel processing** to reduce execution time.

Knowledge graphs on a large scale are at the frontier of AI research. Alas, real-world knowledge is not structured neatly into a schema but rather unstructured, messy, and organic.

Sign up for The Variable

By Towards Data Science

Every Thursday, the Variable delivers the very best of Towards Data Science: from hands-on tutorials and cutting-edge research to original features you don't want to miss. [Take a look.](#)

[Get this newsletter](#)

Emails will be sent to ruijiera@usc.edu.
[Not you?](#)

