

```

1 %pip install requests
2 %pip install pandas
3 %pip install networkx
4 import requests
5 import json
6 import os
7 import time
8 import networkx as nx
9 import pdb
10 import matplotlib.pyplot as plt
11 import pickle
12 import textwrap
13
14 from google.colab import files
15 from google.colab import userdata

```

--NORMAL--

```

➞ Requirement already satisfied: requests in /usr/local/lib/python3.10/dist-packages (2.31.0)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests) (3.3.2)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests) (3.6)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.10/dist-packages (from requests) (2.0.7)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests) (2023.11.17)
Requirement already satisfied: pandas in /usr/local/lib/python3.10/dist-packages (1.5.3)
Requirement already satisfied: python-dateutil>=2.8.1 in /usr/local/lib/python3.10/dist-packages (from pandas) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.10/dist-packages (from pandas) (2023.3.post1)
Requirement already satisfied: numpy>=1.21.0 in /usr/local/lib/python3.10/dist-packages (from pandas) (1.23.5)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil>=2.8.1->pandas) (1.16.0)
Requirement already satisfied: networkx in /usr/local/lib/python3.10/dist-packages (3.2.1)

```

```
1 S2_API_KEY = userdata.get('APIKEY')
```

```

1 authorsNames = []
2 articlesNames = []
3 # Used to prevent getting authors from another field with the same name
4 Field = "Computer Science"

```

Upload custom author list

```

1 uploaded = files.upload()
2 name = list(uploaded.keys())[0]
3
4 with open(f'{name}') as f:
5     authorsNames = f.read().splitlines()

```

Choose files articles.txt

- articles.txt(text/plain) - 654 bytes, last modified: 27/11/2023 - 100% done
- Saving articles.txt to articles.txt

Upload custom article list

```

1 uploaded = files.upload()
2 name = list(uploaded.keys())[0]
3
4 with open(f'{name}') as f:
5     articlesNames = f.read().splitlines()

```

Choose files articles.txt

- articles.txt(text/plain) - 654 bytes, last modified: 02/12/2023 - 100% done
- Saving articles.txt to articles (1).txt

```

1 authors = []
2 qualityWorks = []
3 references = []

```

Getting authors Id's.

WARNING: There's no guarantee that the fetched authors will be the desired ones (the code prioritize the most famous ones in the field that match the name), if some data looks strange, checking on the Semantic Scholar website is recommended.

```

1 for author in authorsNames:
2     first,midle,*last = author.split(' ')
3     last = last if last else " "
4
5 # Don't stop until get the authors, ignoring failures
6 while True:
7     try:
8         time.sleep(.3)
9         rsp = requests.get(f"https://api.semanticscholar.org/graph/v1/author/search?query={first}+{midle}+{last}",
10                             headers={'X-API-KEY': S2_API_KEY},
11                             params={'fields': 'authorId,name,citationCount,papers,papers.fieldsOfStudy','limit': '50'})
12
13         rsp.raise_for_status()
14         break
15
16 except requests.exceptions.RequestException:
17     pass
18
19
20 # add the first author on the search results
21 data = rsp.json()['data']
22 data = sorted(data, key=lambda x:x['citationCount'])
23 data.reverse()
24
25 # get correct author by field of study
26 i = 0
27
28 if not data:
29     break
30
31 while True:
32
33     if data and data[i] and data[i]['papers']:
34         # Try to find a paper that satisfies the conditions
35         iteration_count, found_paper = next(
36             (
37                 (index, paper)
38                 for index, paper in enumerate(data[i]['papers'])
39                 if (paper['fieldsOfStudy'] is not None) and (Field in paper['fieldsOfStudy'])
40             ),
41             (None, None)
42         )
43
44         # Check if a valid paper was found or if the maximum iterations are reached
45         if found_paper is not None and iteration_count < 10:
46             #print(f"Iteration {iteration_count}: {found_paper}")
47             print(f"{data[i]['name']}\t citations: {data[i]['citationCount']} papers: {len(data[i]['papers'])}")
48
49             del data[i]['papers']
50             authors.append(data[i])
51             break
52
53 # Move to the next index
54 if i < len(data)-1:
55     i += 1
56 # If failed to find go to the next on the list
57 else:
58     break
59
60

```

Getting the articles

```

1 for article in articlesNames:
2     try:
3         time.sleep(0.3)
4         rsp = requests.get(f"https://api.semanticscholar.org/graph/v1/paper/autocomplete?query={article}",
5                             headers={'X-API-KEY': S2_API_KEY},
6                             params={'limit': '5'})
7
8         rsp.raise_for_status()
9         if rsp.json()["matches"]:
10
11             item = rsp.json()['matches'][0]
12
13             rsp = requests.get(f"https://api.semanticscholar.org/graph/v1/paper/{item['id']}",
14                                 headers={'X-API-KEY': S2_API_KEY},
15                                 params={'fields': 'title,paperId,fieldsOfStudy,references,references.title,references.citationCount,references.ir
16             rsp.raise_for_status()
17
18             paper = rsp.json()
19
20             if paper['paperId'] is not None:
21                 # add all papers cited by the paper
22                 ref = []
23                 for x in paper['references']:
24                     if x['paperId'] is not None:
25                         x['quality'] = False
26                         ref.append(x)
27
28                 references.append(ref)
29
30                 # add paper
31                 del paper['references']
32                 paper['quality'] = True
33                 qualityWorks.append(paper)
34
35             else:
36                 print(f"Could not find {article}")
37
38 except requests.exceptions.RequestException:
39     pass

```

Getting author's publications and references

```

1 for author in authors:
2     # Get a list of all authors publications and their respective references
3
4     while True:
5         try:
6             time.sleep(0.3)
7             rsp = requests.get(f"https://api.semanticscholar.org/graph/v1/author/{author['authorId']}/papers",
8                                 headers={'X-API-KEY': S2_API_KEY},
9                                 params={'fields': 'title,paperId,fieldsOfStudy,references,references.title,references.citationCount,references.in
10                                     'limit': 900})
11             # reference.citationCount, etc
12
13             rsp.raise_for_status()
14             break
15
16 except requests.exceptions.RequestException:
17     pass
18
19 data = rsp.json()['data']
20
21 for paper in data:
22     if paper['paperId'] is not None:
23         # add all papers cited by the paper
24         ref = []
25         for x in paper['references']:
26             if x['paperId'] is not None:
27                 x['quality'] = False
28                 ref.append(x)
29
30         #print(ref)
31         references.append(ref)
32
33         # add paper
34         del paper['references']
35         paper['quality'] = True
36         qualityWorks.append(paper)
37
38

```

```

1 for paper in qualityWorks:
2     if paper['fieldsOfStudy']:
3         if Field not in paper['fieldsOfStudy']:
4             #print(f"Deleted not in the field {paper['title']}")
5             del paper

```

We now have a list filled with papers and another list filled with a list of that paper references. Thus, basically a adjacency list.

As for now, this list only contains the most important work (the writings of handpicked authors) that we will consider to be of high quality.

Adding depth means to add the references's references into our graph.

```

1 nodes = qualityWorks.copy()
2 edges = references.copy()
3 qualityIds = [x['paperId'] for x in qualityWorks]

```

```

1 # add one more depth into the search
2 def addDepth():
3     lenght = len(edges)
4
5     for i in range(lenght):
6         print(f"{i} / {lenght}")
7         # if has reference
8         if edges[i]:
9             for paper in edges[i]:
10                 end = False
11
12                 # add reference paper into the node list
13                 if paper not in nodes:
14                     nodes.append(paper)
15                 # if is already prossed, go to the next
16                 else:
17                     continue
18
19                 # get paper references
20                 while True:
21                     try:
22
23                         time.sleep(0.1)
24                         rsp = requests.get(f"https://api.semanticscholar.org/graph/v1/paper/{paper['paperId']}/references",
25                                           headers={'X-API-KEY': S2_API_KEY},
26                                           params={'fields': 'title,fieldsOfStudy,citationCount,influentialCitationCount'})
27                         rsp.raise_for_status()
28                         break
29
30                     except requests.exceptions.HTTPError as err:
31                         if rsp.status_code == 404:
32                             end = True
33
34                     except requests.exceptions.RequestException as e:
35                         print(e)
36                         pass
37
38                 if end:
39                     break
40
41                 data = rsp.json()['data']
42                 #print(paper['title'],len(data))
43
44                 # add the paper's references as his edges
45                 ref = []
46                 for x in data:
47                     if x['citedPaper']['paperId'] is not None and x['citedPaper']['fieldsOfStudy']:
48                         if Field in x['citedPaper']['fieldsOfStudy']:
49                             if x['citedPaper']['paperId'] in qualityIds:
50                                 x['citedPaper']['quality'] = True
51                             else:
52                                 x['citedPaper']['quality'] = False
53
54                             ref.append(x['citedPaper'])
55
56                 edges.append(ref)

```

```

1 for i in range(2):
2     addDepth()

```

148 / 204
149 / 204
150 / 204
151 / 204
152 / 204
153 / 204
154 / 204
155 / 204
156 / 204
157 / 204
158 / 204
159 / 204
160 / 204
161 / 204
162 / 204
163 / 204
164 / 204
165 / 204
166 / 204
167 / 204
168 / 204
169 / 204
170 / 204
171 / 204
172 / 204
173 / 204
174 / 204
175 / 204
176 / 204
177 / 204
178 / 204
179 / 204
180 / 204
181 / 204
182 / 204
183 / 204
184 / 204
185 / 204
186 / 204
187 / 204
188 / 204
189 / 204
190 / 204
191 / 204
192 / 204
193 / 204
194 / 204
195 / 204
196 / 204
197 / 204
198 / 204
199 / 204
200 / 204
201 / 204
202 / 204
203 / 204

```
1 print(len(nodes))  
2 print(len(edges))
```

1927
1927

Create or upload(next cell) a graph

```

1 G = nx.DiGraph()
2
3 atributes = ['citationCount','influentialCitationCount','quality','title']
4
5 for i,l in enumerate(edges):
6     for edge in l:
7
8         if edge['paperId'] in qualityIds:
9             edge['quality'] = True
10
11     # Calculate the edge weight
12     if nodes[i]['quality'] == True:
13         influence = 0
14     else:
15         # Reversing the weights so that it has a inverse effect
16         if nodes[i]['influentialCitationCount'] and nodes[i]['citationCount']:
17             influence = 1/(nodes[i]['citationCount'] + 5*nodes[i]['influentialCitationCount'])
18         else:
19             if nodes[i]['citationCount']:
20                 influence = 1/(nodes[i]['citationCount'])
21             else:
22                 influence = 0.5
23
24     ne = {'influence': influence}
25
26     # Add references as edges and nodes, as needed
27     if edge not in nodes:
28         nn = {key: edge[key] for key in atributes}
29
30         if edge['paperId'] in qualityIds:
31             nn['quality'] = True
32
33         G.add_nodes_from([(edge['paperId'],nn)])
34         G.add_edges_from([(edge['paperId'],nodes[i]['paperId'],ne)])
35
36     else:
37         G.add_edges_from([(edge['paperId'],nodes[i]['paperId'],ne)])
38
39 # Add the rest of nodes
40 for node in nodes:
41     ne = {key: node[key] for key in atributes}
42     G.add_nodes_from([(node['paperId'],ne)])
43
44 with open('literature.gpickle', 'wb') as f:
45     pickle.dump(G, f, pickle.HIGHEST_PROTOCOL)
46
47 files.download("literature.gpickle")

```

Upload a graph

```

1 uploaded = files.upload()
2 name = list(uploaded.keys())[0]
3
4 with open(name, 'rb') as f:
5     G = pickle.load(f)

```

Paper ranking

```

1 ranking = []
2
3 for node,data in G.nodes(data=True):
4     # skip calculation on quality papers
5     if data['quality']:
6         continue
7
8     sum = 0
9     counter = 0
10    for work in qualityWorks:
11
12        # calculate distance using dijkstra
13        try:
14            sum = sum + nx.shortest_path_length(G, source=node, target=work['paperId'], weight='influence')
15            # base number is lower when there's no path
16            # except on MutualNeighbors
17
18    ranking = sorted(ranking, key=lambda x:x['score'])
19    ranking.reverse()
20
21    3
22    4 top = ranking[:20]
23    5
24    6 topIds = [x['id'] for x in top]
25    7
26    8 subgraph = G.subgraph(qualityIds + topIds)
27    9
28    10 # quality -> red
29    11 # influent (top 10) -> orange
30    12 # the rest -> blue
31    13 node_colors = ['orange' if node in topIds else 'red' if node in qualityIds else 'blue' for node,data in subgraph.nodes(data=True)]
32    14 score = 1
33
34    1 for paper in top:
35    2     print(paper['title'])

```

ImageNet classification with deep convolutional neural networks
 Very Deep Convolutional Networks for Large-Scale Image Recognition
 Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks
 Learning Multiple Layers of Features from Tiny Images
 Microsoft COCO: Common Objects in Context
 Going deeper with convolutions
 ImageNet Large Scale Visual Recognition Challenge
 Fully convolutional networks for semantic segmentation
 Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift
 Random Forests
 Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation
 Fast R-CNN
 Auto-Encoding Variational Bayes
 ImageNet: A large-scale hierarchical image database
 The Pascal Visual Object Classes (VOC) Challenge
 Caffe: Convolutional Architecture for Fast Feature Embedding
 Rectified Linear Units Improve Restricted Boltzmann Machines
 Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification
 Understanding the difficulty of training deep feedforward neural networks
 U-Net: Convolutional Networks for Biomedical Image Segmentation

```

1 # Draw the graph
2 pos = nx.spring_layout(subgraph, k=1.7)
3 nx.draw(subgraph, pos, node_color=node_colors, font_color="black", font_weight="bold", edge_color="gray", linewidths=1)
4
5 #node_labels = {node: '\n'.join(textwrap.wrap(data['title'], width=26)) for node,data in subgraph.nodes(data=True)}
6 #nx.draw_networkx_labels(G, pos, labels=node_labels,font_size=6, verticalalignment="bottom")
7
8 plt.show()

```

