

산업인공지능 개론

Home Work #2 *Knowledge Graph*

학과 : 산업인공지능학과

학번 : 2024254022

이름 : 정현일

2024.03.31.

KnowledgeGraph.py (1/2)

```
1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3  """
4  산업인지능 개론
5
6  HW2
7  강의노트 3.2장에서 Knowledge Graph를 구성하는 코드를 실행하고, 구성된 knowledge graph에서
8  written_by나 composed_by가 아닌 2개의 관계를 선택하여 해당되는 정보를 추출
9
10 학번 : 2024254022
11 이름 : 정현일
12
13 @author: chohi
14 """
15
16
17 import re
18 import pandas as pd
19 import bs4
20 import requests
21 import spacy
22
23 from spacy import displacy
24 nlp = spacy.load('en_core_web_sm')
25
26 from spacy.matcher import Matcher
27 from spacy.tokens import Span
28
29 import networkx as nx
30
31 import matplotlib.pyplot as plt
32 from tqdm import tqdm
33
34
35 pd.set_option('display.max_colwidth', 200)
36 %matplotlib inline
37
38 candidate_sentences = pd.read_csv('wiki_sentences_v2.csv')
39 print(candidate_sentences.shape)
40 print(candidate_sentences)
41
42
43 entity_pairs = []
44
45
46 def get_entities(sent):
47     ent1 = ""
48     ent2 = ""
49     prv_tok_dep = ""
50     prv_tok_text = ""
```

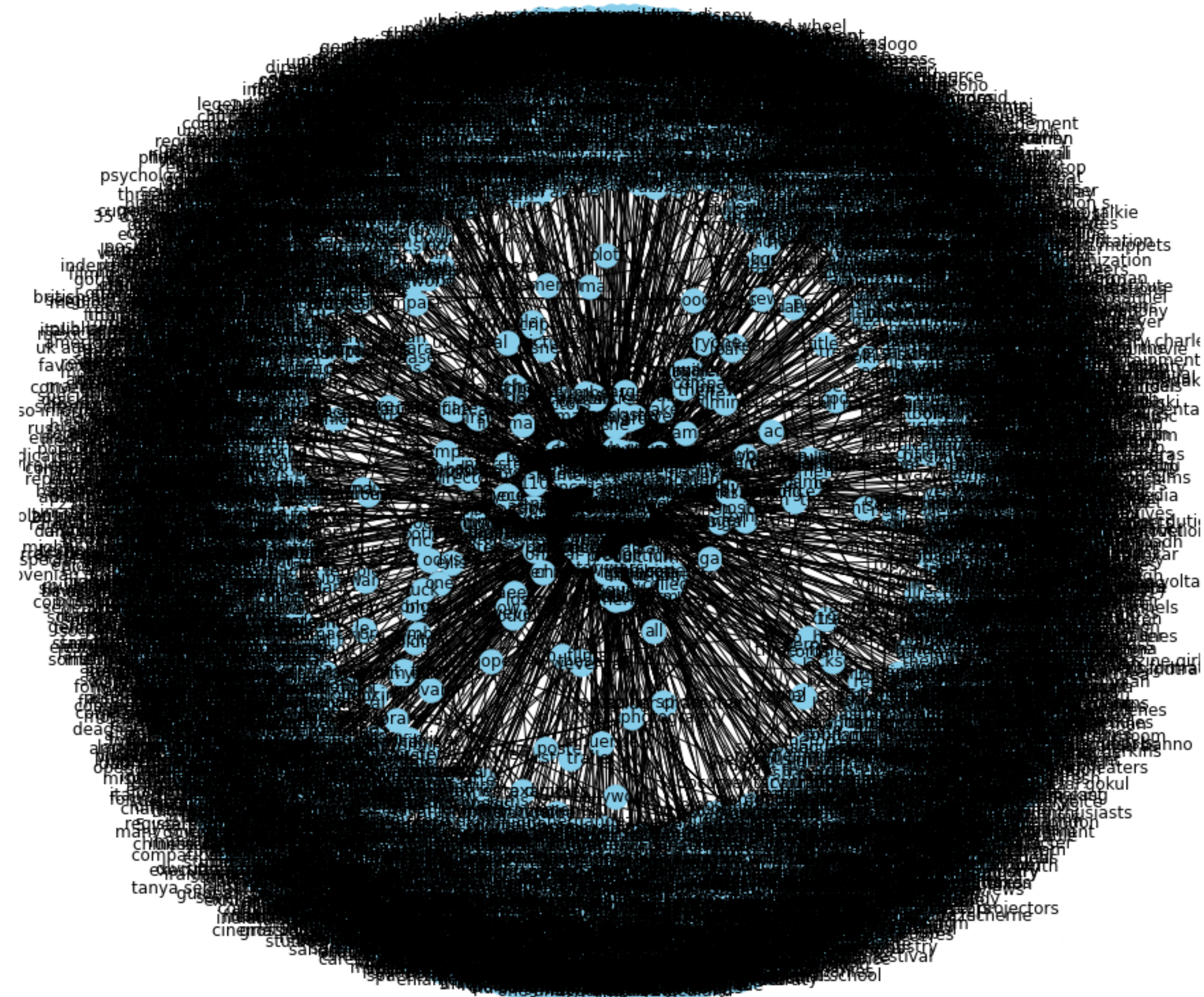
```
51     prv_tok_text = ""
52     prefix = ""
53     modifier = ""
54     for tok in nlp(sent):
55         # 토큰이 구두점(punctuation mark)이면 다음 토큰으로 이동
56         if tok.dep_ != "punct":
57             if tok.dep_ == "compound": # 토큰이 복합어인 경우
58                 prefix = tok.text
59                 if prv_tok_dep == "compound": # 직전 토큰이 복합어이면 현재 토큰과 결합
60                     prefix = prv_tok_text + " " + tok.text
61             if tok.dep_.endswith("mod") == True: # 토큰이 수식어(modifier)인 경우
62                 modifier = tok.text
63                 if prv_tok_dep == "compound": # 직전 토큰이 수식어이면 현재 토큰을 결합
64                     modifier = prv_tok_text + " " + tok.text
65
66             if tok.dep_.find("subj") == True: # 주어(subject)인 경우,
67                 ent1 = modifier + " " + prefix + " " + tok.text # 수식어와 현재 토큰 결합 => 개체명 생성
68                 prefix = ""
69                 modifier = ""
70                 prv_tok_dep = ""
71                 prv_tok_text = ""
72             if tok.dep_.find("obj") == True: # 목적어인 경우
73                 ent2 = modifier + " " + prefix + " " + tok.text # 수식어와 현재 토큰 결합 => 객체명 생성
74
75             prv_tok_dep = tok.dep_
76             prv_tok_text = tok.text
77     return [ent1.strip(), ent2.strip()] # 식별된 개체명 반환
78
79
80
81 def get_relation(sent):
82     doc = nlp(sent)
83     matcher = Matcher(nlp.vocab)
84
85     # 패턴 정의
86     pattern = [{ 'DEP' : 'ROOT'},
87                 { 'DEP' : 'prep', 'OP' : "?"},
88                 { 'DEP' : 'agnet', 'OP' : "?"},
89                 { 'POS' : 'ADJ', 'OP' : "?"}]
90
91     matcher.add('matching_1', [pattern])
92
93     matches = matcher(doc)
94     print('matches : ', matches)
95     k = len(matches) - 1
96
97     span = doc[matches[k][1]:matches[k][2]]
98
99     return(span.text)
100
101
```

KnowledgeGraph.py (2/2)

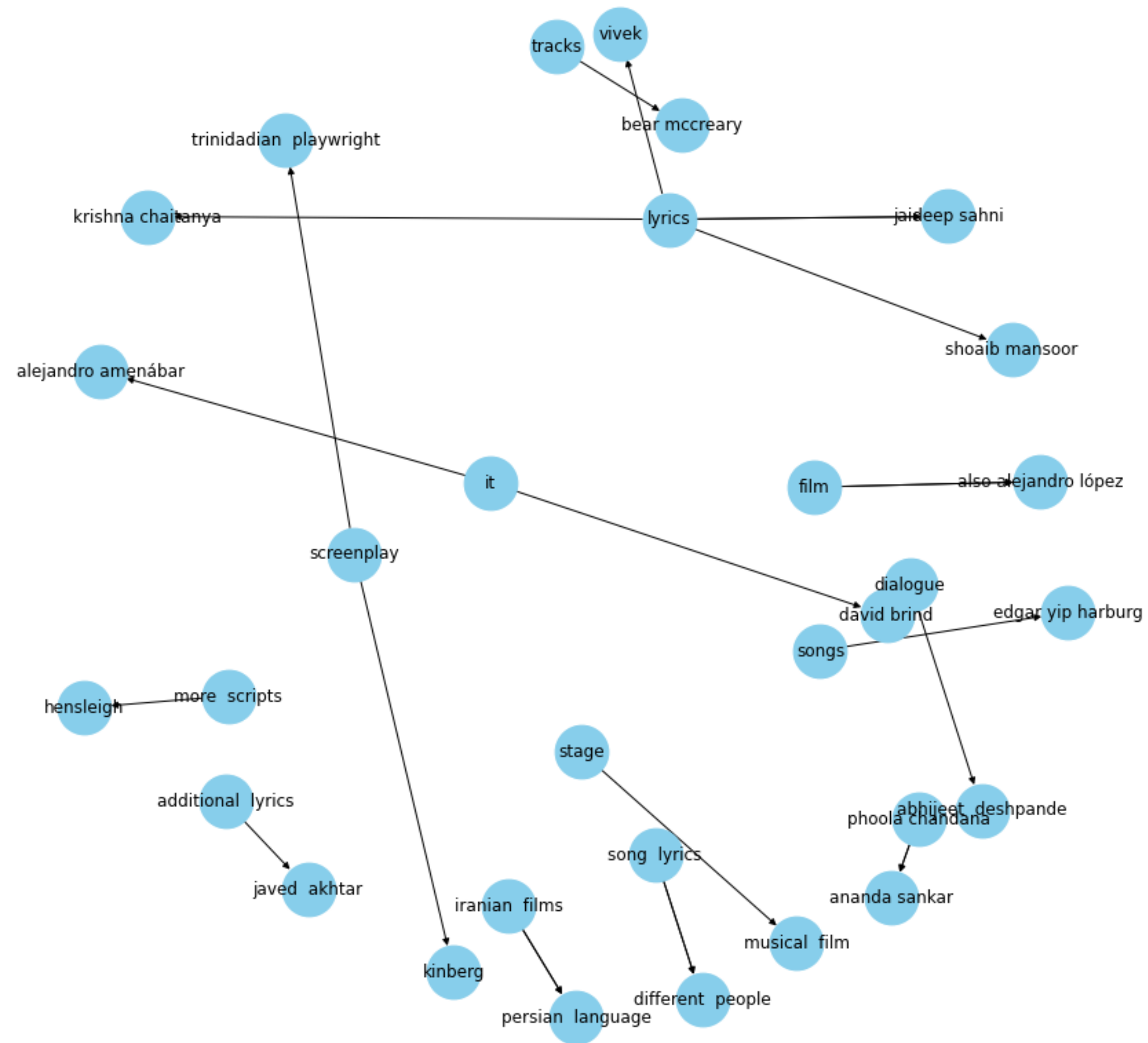
```
100
101
102 get_relation("John completed the task")
103 relations = [get_relation(i) for i in tqdm(candidate_sentences['sentence'])]
104
105
106
107 for i in tqdm(candidate_sentences["sentence"]):
108     entiry_pairs.append(get_entities(i))
109
110 entiry_pairs[10:20]
111
112
113 # 주어(subject) 추출
114 source = [i[0] for i in entiry_pairs]
115
116
117 # 목적어(object) 추출
118 target = [i[1] for i in entiry_pairs]
119
120 kg_df = pd.DataFrame({'source' : source, 'target' : target, 'edge' : relations})
121
122
123
124 # 방향 그래프 생성
125 G = nx.from_pandas_edgelist(kg_df, "source", "target",
126                             edge_attr=True, create_using=nx.MultiDiGraph())
127
128 # 그래프 그리기
129 plt.figure(figsize=(12,12))
130 pos = nx.spring_layout(G)
131 nx.draw(G, with_labels=True, node_color='skyblue', edge_cmap=plt.cm.Blues, pos=pos)
132 plt.show()
133
134
135
136 ## written
137 G=nx.from_pandas_edgelist(kg_df[kg_df['edge']=="written"], "source", "target",
138                             edge_attr=True, create_using=nx.MultiDiGraph())
139
140 plt.figure(figsize=(12,12))
141 pos = nx.spring_layout(G, k = 0.5)
142 nx.draw(G, with_labels=True, node_color='skyblue', node_size=1500, edge_cmap=plt.cm.Blues, pos = pos)
143 plt.show()
144
145
146
147 ## include
148 B = nx.from_pandas_edgelist(kg_df[kg_df['edge'] == "include"], "source", "target",
149                             edge_attr=True, create_using=nx.MultiDiGraph())
150
151
```

```
151 plt.figure(figsize=(12,12))
152 pos = nx.spring_layout(G, k=0.5) # k refulates the distance between nodes
153 nx.draw(G, with_labels=True, node_color='skyblue', node_size=1500, edge_cmap=plt.cm.Blues, pos = pos)
154 plt.show()
155
```


실행 결과



실행 결과



실행 결과

