

PROJECT PRESENTATION

Created by Francesco Villi

TASK

From a dataset of movies extract all data, create a graph using Networkx, and resolve the exercises written in the next slides.

Before starting we need to create the graph...

CREATION GRAPH

CREATION GRAPH

- Split values actor and movie
- Check if the file is ended
- Detect film, year and actor
- Create or retrieve the key for film and actor

```
1 def createGraph():
2     EOF=False
3     id_actor=0
4     id_film=1
5     count_actor=0
6     count_film=1
7     actor_to_key={}
8     key_to_actor={}
9     film_to_key={}
10    key_to_film={}
11    G=nx.Graph()
12    while not EOF:
13        line=file.readline().split('\n')
14        record=line[0].split('\t')
15        if len(record)<2:
```

RESULT

- The graph has 3.110.737 nodes where 2.364.796 are actors and 745.941 are movies
- The amount of edges is 8103960
- The algorithm showed in the previous slide takes about 1 minute to be executed

EXERCISE 1

Which is the movie with the largest number of actors,
considering only the movies up to year $x =$
 $\{1930, 1940, 1950, 1960, 1970, 1980, 1990, 2000, 2010, 2020\}$?

ALGORITHM

- For each film, check the year
- Calculate the number of actors
- Update the best result

```
1 x = [1930,1940,1950,1960,1970,1980,1990,2000,2010,2020]
2 up_to_year=x[rnd.randint(0,len(x)-1)]
3 best_result={'count':0,'movies':[]}
4 for film,name_year in key_to_film.items():
5     if name_year[1]<=up_to_year:
6         tot_movie=len(G[film])
7         if(tot_movie>=best_result['count']):
8             if tot_movie>best_result['count']:
9                 best_result['movies']=[]
10                best_result['count']=tot_movie
11                best_result['movies'].append(film)
```

RESULTS

Year	Film	Nr actors
1930	The King of Kings (1927)	171
1940	The Buccaneer (1938)	219
1950	Gone to Earth (1950)	290
1960	Around the World in Eighty Days (1956)	1298
1970	Around the World in Eighty Days (1956)	1298
1980	Around the World in Eighty Days (1956)	1298
1990	Around the World in Eighty Days (1956)	1298
2000	Around the World in Eighty Days (1956)	1298
2010	Around the World in Eighty Days (1956)	1298
2020	Around the World in Eighty Days (1956)	1298

EXERCISE 2

Considering only the movies up to year $x = \{1930, 1940, 1950, 1960, 1970, 1980, 1990, 2000, 2010, 2020\}$ and restricting to the largest connected component of the graph.
Compute exactly the diameter of G

This Algorithm is based on 3 steps:

- Method for finding the largest connected component H according to the film's year
- two_sweep for detecting the starting node in the subgraph H
- iFUB for computing the exact diameter of H

TWO_SWEEP

```
1 @py_random_state(1)
2 def two_sweep(Graph, seed):
3     rnd_node = seed.choice(list(Graph))
4     source= list(nx.single_source_shortest_path_length(Graph, rnd_node))[-1]
5     distances_b = list(nx.single_source_shortest_path_length(Graph, source).items())
6     index_start_node=bisect.bisect_left(distances_b,int(distances_b[-1][1]/2),key=lambda k:k[1])
7     index_end_node=bisect.bisect_right(distances_b,int(distances_b[-1][1]/2),key=lambda k:k[1])
8     return max(distances_b[index_start_node:index_end_node],key=lambda k:len(H_small[k[0]]))
```

- Select a random source node and get the farthest node from the source
- Compute the shortest path lengths from source to all reachable nodes.
- Select a node in the middle with the highest degree

iFUB

```
1 def iFUB(G,node_start):
2     i=nx.eccentricity(G,node_start)
3     lb=i
4     ub=2*i
5     list_fringe=itertools.groupby(reversed(nx.single_source_shortest_path_length(
6         G,node_start,cutoff=i).items()),key=lambda k:k[1])
7     max_ecc=0
8     while ub>lb:
```

- Compute the eccentricity of the starting node and set a lower and upper bound
- Get a list of nodes until distance i from the starting node
- For each node at distance i, compute eccentricity and check if it's equal or greater than upper bound
- If the previous condition isn't triggered, update the lower and upper bound

ALGORITHM

```
1 x = [1930,1940,1950,1960,1970,1980,1990,2000,2010,2020]
2 year=x[rnd.sample(x,1)]
3 list_nodes=set()
4 for u,v in G.edges():
5     if G[u][v]['year'] <= year:
6         list_nodes.add(u)
7         list_nodes.add(v)
8 largest_cc=max(nx.connected_components(G.subgraph(list_nodes)),key=len)
9 H_small=G.subgraph(largest_cc)
10 node_start=two_sweep(H_small,None)
11 iFUB(H_small,node_start[0])
```

- Find all nodes that respect the year's constraint
- Find maximum connected component and create a SubGraph view of the subgraph induced on nodes
- Compute the diameter of the subgraph

Observations

- In some cases, the computing of the diameter can take much time due to the unlucky starting node

RESULTS

Year	Diameter	Time
1930	32	50 sec
1940	36	54 sec
1950	38	92 sec
1960	34	104 sec
1970	25	11 min
1980	25	31 min
1990	30	4 min
2000	28	13 min
2010	30	68 min
2020	32	44 min

EXERCISE 3

**Which is the movie with the largest number of popular actors,
i.e. such that the sum of the number of movies its actors
participated in is maximum?**

ALGORITHM

- For each film retrieve the list of actors
- For each actor in that list, compute the popularity
- Compare each sum with the best result

```
1 largest_movie={'count':0,'movies':[]}
2 for film,_ in key_to_film:
3     edges=G.edges(film)
4     count=sum(G.degree(act) for flm,act in G.edges(film))
5     if count>=largest_movie['count']:
6         if count>largest_movie['count']:
7             largest_movie['count']=count
8             largest_movie['movies']=[]
9             largest_movie['movies'].append(film)
10 print(largest_movie)
```

RESULTS

Film	Popularity
MILF Madness (2012) (V)	34181

EXERCISE 4

Build also the actor graph, whose nodes are only actors and two actors are connected if they did a movie together, which is the pair of actors who collaborated the most among themselves?

CREATE ACTOR GRAPH

```
1 def create_graph_actors_ordered():
2     collaborations={}
3     best_collaborators=[0,[]]
4     G_actor=nx.Graph()
5     for actor_1 in key_to_actor:
6         list_films=itertools.chain(G[actor_1])
7         for film in list_films:
8             list_actor_for_film=itertools.chain(G[film])
9             for actor_2 in list_actor_for_film:
10                 if actor_2>actor_1: # avoid duplication
11                     if actor_2 in collaborations:
```

- For each actor get the list of films in which he participated
- For each film, get the list of actor in that film
- Compute the collaborations between actors maintaining the greatest
- Add others and best collaborations to the graph

ALGORITHM

```
1 best_result={'count':0,'pairs':[]}
2 for actor in itertools.chain(G_actor):
3     for actor_2 in itertools.chain(reversed(list(G_actor[actor]))):
4         colls=G_actor[actor][actor_2]['weight']
5         if colls>=best_result['count']:
6             if colls>best_result['count']:
7                 best_result['count']=colls
8                 best_result['pairs']=[]
```

- For each actor, get the reversed list of adjacent nodes
- Check the weight for each edge that links the two actors
- We start to compare the best collaboration of that actor with the best result
- If the count of collaborators is lower than the best result, we skip to the next actor

RESULTS

Actors	Collaborations	Algorithm	Time
Byron, Tom (I) - North, Peter (I)	420	create_graph_actors_ordered	6 min
		find_best_collaborations	4 min