

✓ Homework Assignment 2

[The Art of Analyzing Big Data - The Data Scientist's Toolbox](#)

By Dr. Michael Fire

✓ Dataset Collecting

Question 1: Write a function that collects all titles and number of votes for each title of a given [hacker news page](#) (15pt)

```
1 import requests
2
3 def get_hacker_news_titles_and_votes():
4     url = 'https://hacker-news.firebaseio.com/v0/topstories.json?print=pretty'
5     response = requests.get(url)
6
7     if response.status_code == 200:
8         top_story_ids = response.json()[:20] # Fetching top 20 stories for demonstration
9         titles_and_votes = []
10
11         for story_id in top_story_ids:
12             story_url = f'https://hacker-news.firebaseio.com/v0/item/{story_id}.json?print=pretty'
13             story_response = requests.get(story_url)
14
15             if story_response.status_code == 200:
16                 story_data = story_response.json()
17                 title = story_data.get('title', '')
18                 votes = story_data.get('score', 0)
19                 user = story_data.get('by', '')
20                 age = story_data.get('time', '')
21
22                 titles_and_votes.append({
23                     'title': title,
24                     'votes': votes,
25                     'user': user,
26                     'age': age
27                 })
28
29         return titles_and_votes
30     else:
31         print(f"Error fetching top stories. Status code: {response.status_code}")
32         return []
33
34 # Example usage
35 hacker_news_data = get_hacker_news_titles_and_votes()
36
37 for index, item in enumerate(hacker_news_data, start=1):
38     print(f"{index}. {item['title']} - {item['votes']} points by {item['user']} {item['age']}")
39
```

1. Show HN: #!/usr/bin/env docker run - 331 points by adtac 1705202405
2. Citadel, a Calibre-compatible eBook management app - 177 points by phildenhoff 1705213360
3. Type information for faster Python C extensions - 43 points by ingve 1705225728
4. What Was ISDN? - 96 points by ecliptik 1705121345
5. OpenD, a D language fork that is open to your contributions - 174 points by mepian 1705196099
6. When "blocked indefinitely" is not indefinite - 17 points by g0xA52A2A 1705228885
7. The Ultimate Docker Cheat Sheet - 74 points by jpeer264 1705225524
8. Smart binoculars can identify 9k birds - 119 points by thunderbong 1705040052
9. BuildZoom (better way to build custom homes) Is hiring a Growth Associate - 1 points by the_economist 1705233860
10. Towards Modern Development of Cloud Applications (2023) - 96 points by signal1 1705216656
11. Building a fully local LLM voice assistant to control my smart home - 534 points by JohnTheNerd 1705183930
12. Posthog is closing their Slack community in favor of forum - 232 points by vmatsiiako 1705205576
13. Dynamic Programming is not Black Magic - 56 points by qsantos 1705225335
14. Ask HN: Does Cloudflare block HN comments if you have code blocks in a reply? - 245 points by zikduruqe 1705192250
15. Thinking in an array language - 238 points by tosh 1705164209
16. Looking back at Postgres (2019) - 14 points by infra_dev 1705038205
17. John Michell: Country Parson Described Black Holes in 1783 (2000) - 128 points by mikequinlan 1705189863
18. Nuclear battery produces power for 50 years without needing to charge - 26 points by prakhar897 1705233993
19. Free unexpected MIT courses to kick start the new year - 73 points by Anon84 1705225776
20. Stellarium: Software which renders realistic skies in real time - 660 points by tosh 1705162365

Question 2: Write a function that collects data on four Marvel movies from [The Movie Database](#). You can use [tmdbv3api](#) (15pt)

```
1 pip install tmdbv3api
```

```
Collecting tmdbv3api
  Downloading tmdbv3api-1.9.0-py3-none-any.whl (25 kB)
Requirement already satisfied: requests in /usr/local/lib/python3.10/dist-packages (from tmdbv3api) (2.31.0)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests->tmdbv3api) (3.3.2)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests->tmdbv3api) (3.6)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.10/dist-packages (from requests->tmdbv3api) (2.0.7)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests->tmdbv3api) (2023.11.17)
Installing collected packages: tmdbv3api
Successfully installed tmdbv3api-1.9.0
```

```
1 from tmdbv3api import TMDb, Movie
2
3 def collect_marvel_movie_data(api_key):
4     tmdb = TMDb()
5     tmdb.api_key = api_key
6
7     movie = Movie()
8     movie_id =[299536, 24428, 569094, 102382] # Example Marvel movie ID (Avengers: Endgame)
9     for i in range(len(movie_id)):
10         print(movie_id[i])
11         m = movie.details(movie_id[i])
12
13         print("\nMovie Details:")
14         print(f"Title: {m.title}")
15         print(f"Overview: {m.overview}")
16         print(f"Popularity: {m.popularity}")
17 # Example usage
18 api_key = '06497a828f86e228abab0d482bf138d4' # Replace with your actual API key
19 collect_marvel_movie_data(api_key)
20
```

299536

Movie Details:

Title: Avengers: Infinity War

Overview: As the Avengers and their allies have continued to protect the world from threats too large for any one hero to handle, a new danger has emerged from the cosmic shadows: Thanos. A desp

Popularity: 236.242

24428

Movie Details:

Title: The Avengers

Overview: When an unexpected enemy emerges and threatens global safety and security, Nick Fury, director of the international peacekeeping agency known as S.H.I.E.L.D., finds himself in need of

Popularity: 178.92

569094

Movie Details:

Title: Spider-Man: Across the Spider-Verse

Overview: After reuniting with Gwen Stacy, Brooklyn's full-time, friendly neighborhood Spider-Man is catapulted across the Multiverse, where he encounters the Spider Society, a team of Spider-Pe

Popularity: 359.731

102382

Movie Details:

Title: The Amazing Spider-Man 2

Overview: For Peter Parker, life is busy. Between taking out the bad guys as Spider-Man and spending time with the person he loves, Gwen Stacy, high school graduation cannot come quickly enough.

Popularity: 100.049

▼ Kickstarter Projects Dataset

Using the [Kickstarter Projects Dataset](#) and Pandas, please answer one of following questions:

Note: Use ks-projects-201801.csv data

▼ Please answer only **one** of the following questions according to your (ID number + 1) (use the formula **mod 3 +1**)

```
1 # which question to answer - put your ID number and run the code
2 your_id = "207376187"
3 q = (int(your_id) + 1) % 3 + 1
4 print("You need to answer questions %s and 4" % q)
```

You need to answer questions 1 and 4

Question 1: On average which project category received the lowest number of backers? (15 pt)

```
1 !pip install kaggle
```

```
1 import pandas as pd
2 from google.colab import files
```

No file chosen Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.

```
Saving ks-projects-201801.csv to ks-projects-201801.csv
cp: cannot stat 'kaggle.json': No such file or directory
chmod: cannot access '/root/.kaggle/kaggle.json': No such file or directory
```

```
Traceback (most recent call last):
  File "/usr/local/bin/kaggle", line 5, in <module>
    from kaggle.cli import main
  File "/usr/local/lib/python3.10/dist-packages/kaggle/__init__.py", line 23, in <module>
    api.authenticate()
  File "/usr/local/lib/python3.10/dist-packages/kaggle/api/kaggle_api_extended.py", line 403, in authenticate
    raise IOError('Could not find {}. Make sure it\'s located in'
OSError: Could not find kaggle.json. Make sure it's located in /root/.kaggle. Or use the environment method.
unzip: cannot find or open kickstarter-projects.zip, kickstarter-projects.zip.zip or kickstarter-projects.zip.ZIP.
```

	ID	name	category	main_category	currency	deadline	goal	launched	pledged	state	backers	country	usd pledged	usd
0	1000002330	The Songs of Adelaide & Abullah	Poetry	Publishing	GBP	2015-10-09	1000.0	2015-08-11 12:12:28	0.0	failed	0	GB	0.0	
1	1000003930	Greeting From Earth: ZGAC Arts Capsule For ET	Narrative Film	Film & Video	USD	2017-11-01	30000.0	2017-09-02 04:43:57	2421.0	failed	15	US	100.0	
2	1000004038	Where is Hank?	Narrative Film	Film & Video	USD	2013-02-26	45000.0	2013-01-12 00:20:50	220.0	failed	3	US	220.0	
ToshiCapital														

```

1 # Group the DataFrame by the 'category' column and calculate the mean number of backers for each category
2 average_backers_by_category = df.groupby('category')['backers'].mean()
3
4 # Find the category with the lowest average number of backers
5 lowest_average_backers_category = average_backers_by_category.idxmin()
6
7 # Display the results
8 print(f"The category with the lowest average number of backers is: {lowest_average_backers_category}")
9 print(f"Average number of backers for this category: {average_backers_by_category[lowest_average_backers_category]}")

```

```

The category with the lowest average number of backers is: Crochet
Average number of backers for this category: 7.851851851851852

```

Question 2: On average which project category received the lowest pledged USD? (15 pt)

```

1

```

Question 3: In which month is there the lowest number of projects? (15 pt)

```

1

```

Question 4 (for all): Visualize your answer using matplotlib or seaborn (15pt)

```

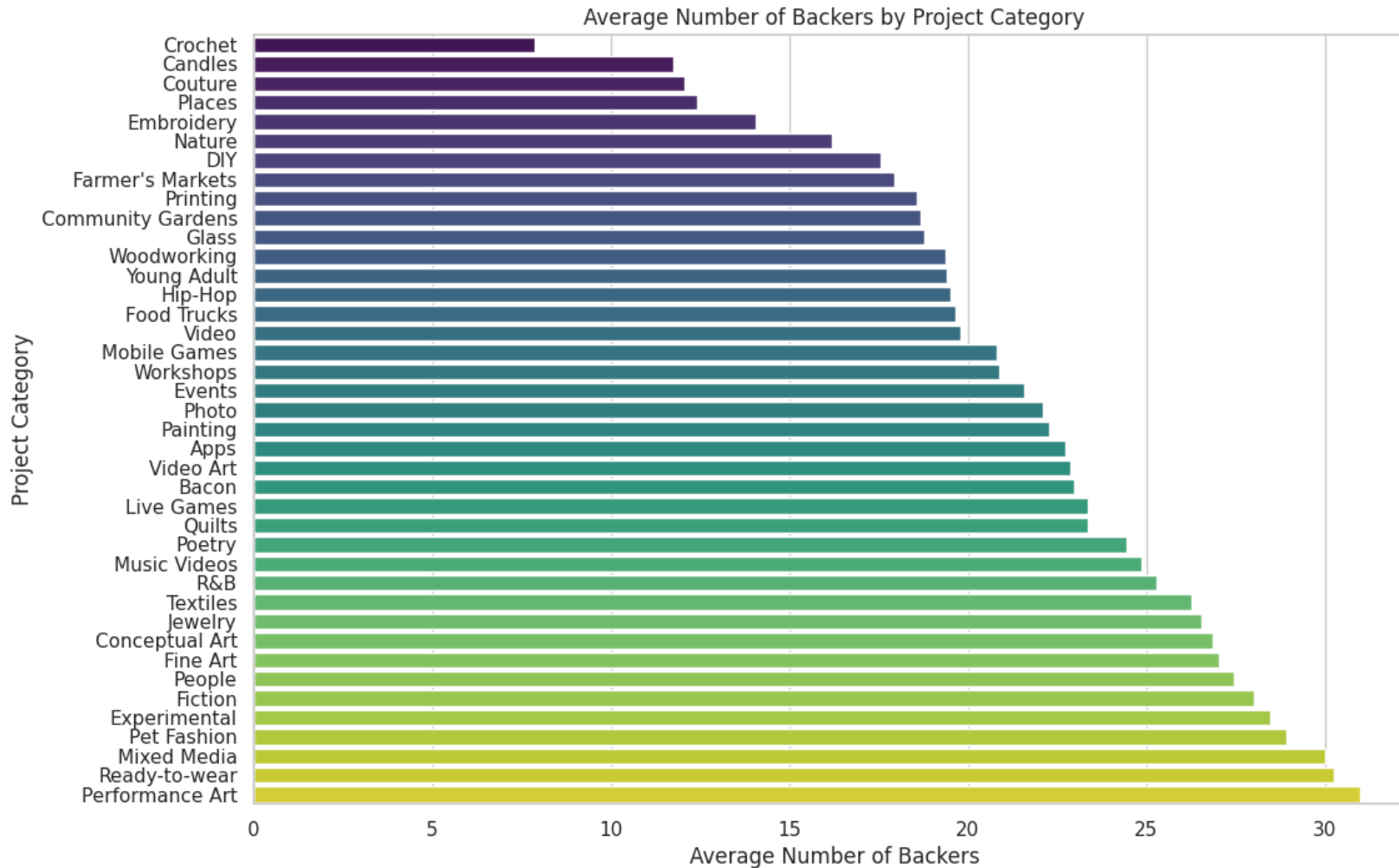
1  import matplotlib.pyplot as plt
2  import seaborn as sns
3
4  # Set the style for the plot
5  sns.set(style="whitegrid")
6
7  # Group the DataFrame by the 'category' column and calculate the mean number of backers for each category
8  average_backers_by_category = df.groupby('category')['backers'].mean().reset_index()
9
10 # Sort the DataFrame by the average number of backers in ascending order
11 average_backers_by_category = average_backers_by_category.sort_values(by='backers', ascending=True)
12
13 # Create a horizontal bar plot
14 plt.figure(figsize=(12, 8))
15 sns.barplot(x='backers', y='category', data=average_backers_by_category.head(40), palette='viridis')
16
17 # Set plot labels and title
18 plt.xlabel('Average Number of Backers')
19 plt.ylabel('Project Category')
20 plt.title('Average Number of Backers by Project Category')
21
22 # Show the plot
23 plt.show()

```

<ipython-input-12-78a13d0dae1e>:15: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `y` variable to `hue` and set `legend=False` for the same effect.

```
sns.barplot(x='backers', y='category', data=average_backers_by_category.head(40), palette='viridis')
```



✓ The Marvel Universe Social Network

Using the [The Marvel Universe Social Network](#) and **Pandas**, please answer the following questions:

Question 1: Write code which calculate the top-20 most friendly characters, i.e., characters with the highest number of friends. Please use *hero_network.csv* file (15pt).

Note: Not all the links in this dataset are symmetric.

```
1 import pandas as pd
2 from google.colab import files
3 # Upload the Kaggle API key
4 files.upload()
5
6 # Move the uploaded Kaggle API key to the required directory
7 !mkdir -p ~/.kaggle
8 !cp kaggle.json ~/.kaggle/
9
10 # Set permissions for the Kaggle API key
11 !chmod 600 ~/.kaggle/kaggle.json
```

Choose Files

No file chosen

Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.

Saving hero-network.csv to hero-network.csv

cp: cannot stat 'kaggle.json': No such file or directory

chmod: cannot access '/root/.kaggle/kaggle.json': No such file or directory

```
1 # Download the dataset from Kaggle
2 !kaggle datasets download -d csanhueza/the-marvel-universe-social-network
```

Traceback (most recent call last):

File "/usr/local/bin/kaggle", line 5, in <module>

from kaggle.cli import main

File "/usr/local/lib/python3.10/dist-packages/kaggle/__init__.py", line 23, in <module>

api.authenticate()

File "/usr/local/lib/python3.10/dist-packages/kaggle/api/kaggle_api_extended.py", line 403, in authenticate

raise IOError('Could not find {}'. Make sure it\'s located in'

OSError: Could not find kaggle.json. Make sure it's located in /root/.kaggle. Or use the environment method.

```
1 # Read the hero network dataset into a Pandas DataFrame
2 hero_network_df = pd.read_csv('hero-network.csv')
3
4 # Display the first few rows of the DataFrame
5 hero_network_df.head()
```

	hero1	hero2
0	LITTLE, ABNER	PRINCESS ZANDA
1	LITTLE, ABNER	BLACK PANTHER/T'CHAL
2	BLACK PANTHER/T'CHAL	PRINCESS ZANDA
3	LITTLE, ABNER	PRINCESS ZANDA
4	LITTLE, ABNER	BLACK PANTHER/T'CHAL

```

1 import pandas as pd
2
3 # Load the hero network dataset
4 hero_network_df = pd.read_csv('hero-network.csv')
5
6 # Drop duplicate connections
7 unique_connections_df = hero_network_df[['hero1', 'hero2']].drop_duplicates()
8
9 # Concatenate 'hero1' and 'hero2' columns
10 all_heroes = pd.concat([unique_connections_df['hero1'], unique_connections_df['hero2']])
11
12 # Count the unique occurrences of each hero
13 unique_connections_count = all_heroes.value_counts()
14
15 # Select the top 20 heroes with the highest number of unique connections
16 top_20_most_connected_characters = unique_connections_count.head(20)
17
18 # Display the result
19 print("Top 20 Most Connected Characters (Unique Friends):")
20 print(top_20_most_connected_characters)

```

Top 20 Most Connected Characters (Unique Friends):

CAPTAIN AMERICA	2854
SPIDER-MAN/PETER PAR	2563
IRON MAN/TONY STARK	2238
WOLVERINE/LOGAN	2036
THING/BENJAMIN J. GR	2023
MR. FANTASTIC/REED R	1958
HUMAN TORCH/JOHNNY S	1939
SCARLET WITCH/WANDA	1932
THOR/DR. DONALD BLAK	1913
BEAST/HENRY & HANK & P	1840
VISION	1814
INVISIBLE WOMAN/SUE	1796
HAWK	1709
CYCLOPS/SCOTT SUMMER	1616
STORM/ORORO MUNROE S	1581
HULK/DR. ROBERT BRUC	1571
ANGEL/WARREN KENNETH	1563
WASP/JANET VAN DYNE	1557
SHE-HULK/JENNIFER WA	1538
DR. STRANGE/STEPHEN	1530

dtype: int64

Question 2: Using Pandas and Networkx create a graph object of The Marvel Universe Social Network with the 150 most "friendly" characters (10pt). The vertices in that graph need to be relative to the size of each character's number of links (also referred to as the vertex degree) (10pt). Please color each node in the graph according to character type according to data in the *nodes.csv* file (5pt)

```

1 # Upload the Kaggle API key
2 files.upload()
3
4 # Move the uploaded Kaggle API key to the required directory
5 !mkdir -p ~/.kaggle
6 !cp kaggle.json ~/.kaggle/
7
8 # Set permissions for the Kaggle API key
9 !chmod 600 ~/.kaggle/kaggle.json

```


[Choose Files](#)

No file chosen

Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.

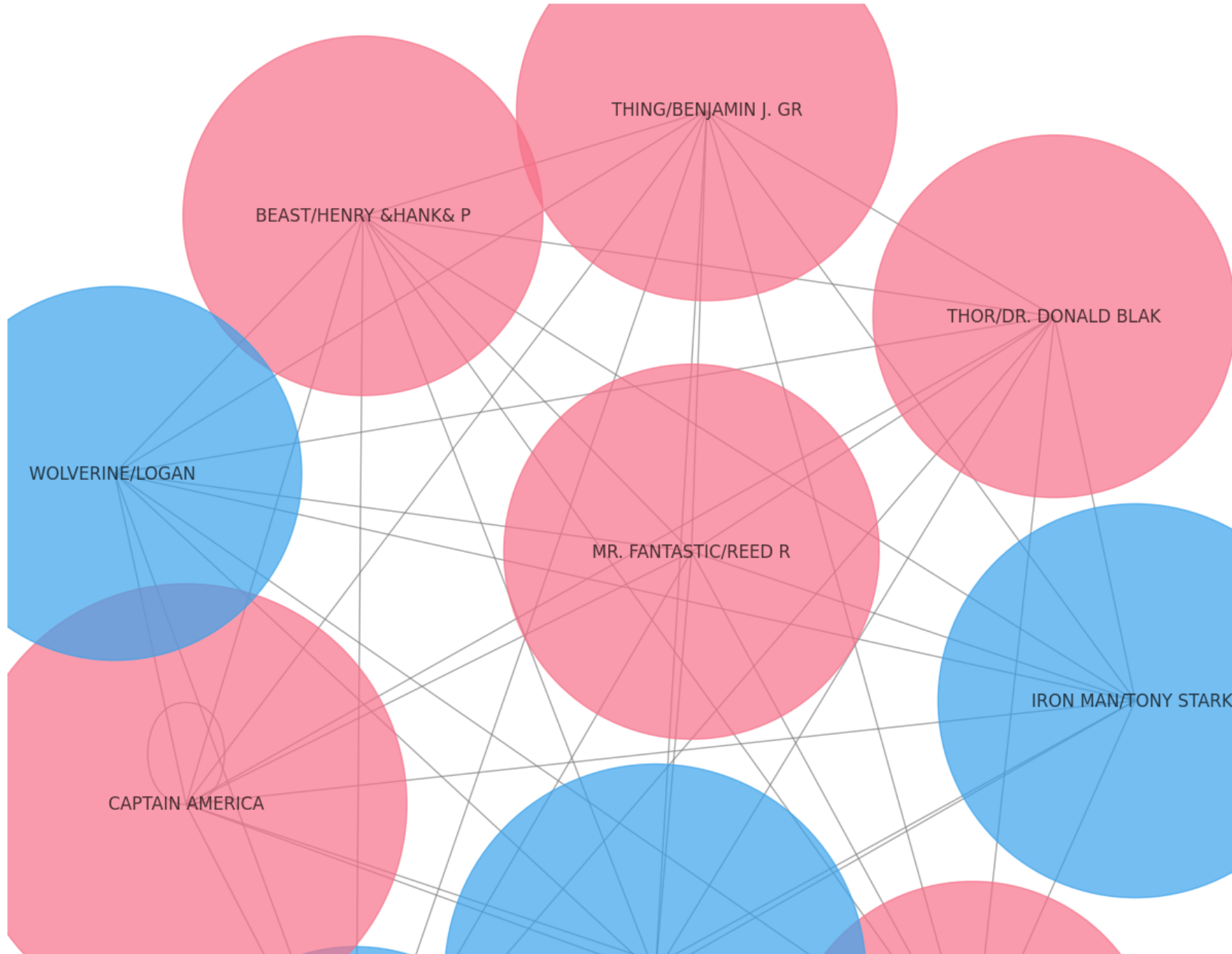
Saving nodes.csv to nodes.csv

cp: cannot stat 'kaggle.json': No such file or directory

chmod: cannot access '/root/.kaggle/kaggle.json': No such file or directory

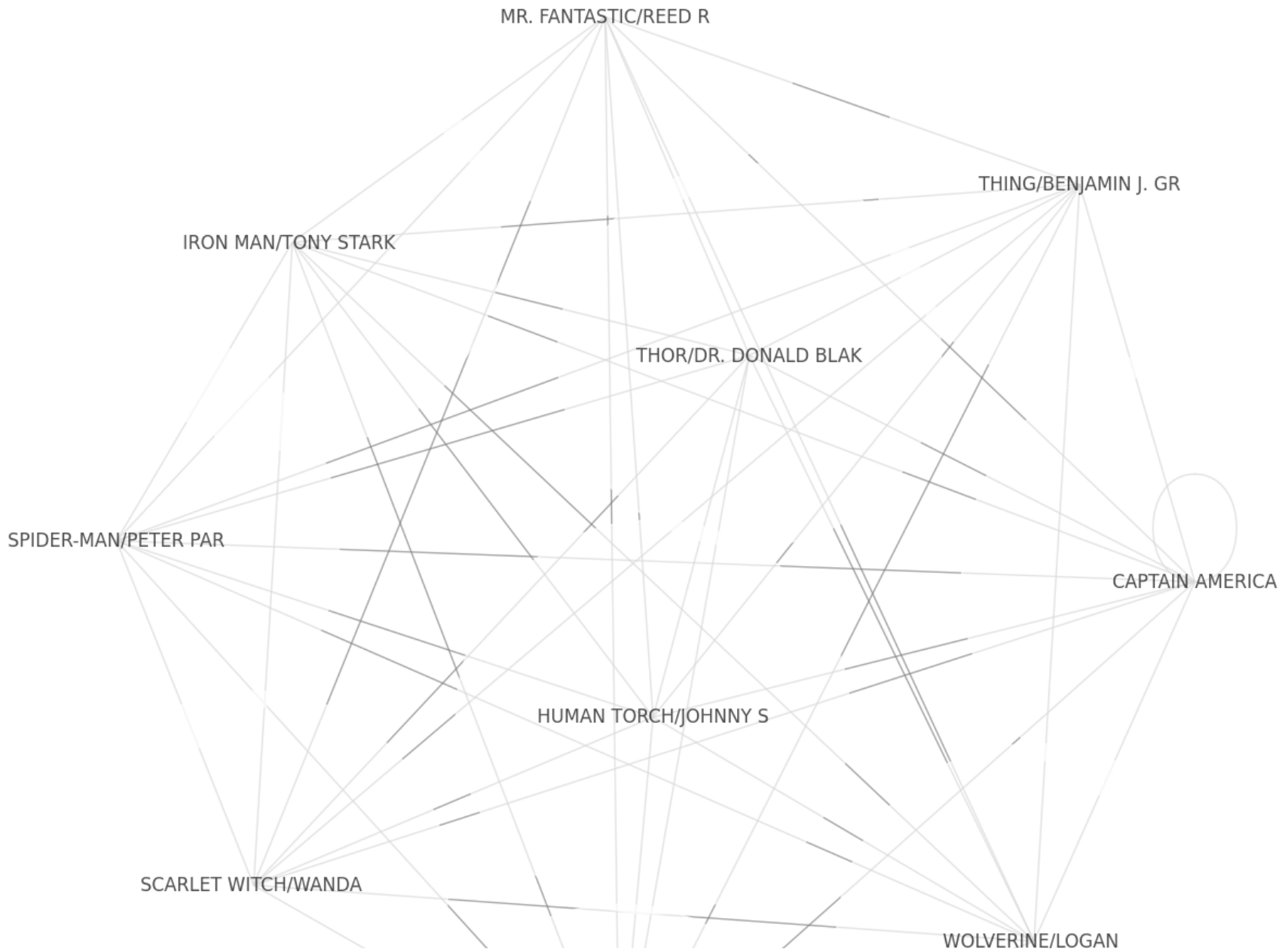
```
1 # Most 10 friendly colored. only for easy demonstraion.
2
3 import pandas as pd
4 import networkx as nx
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7
8 # Load the hero network dataset
9 hero_network_df = pd.read_csv('hero-network.csv')
10
11 # Drop duplicate connections
12 unique_connections_df = hero_network_df[['hero1', 'hero2']].drop_duplicates()
13
14 # Create a graph object
15 G = nx.Graph()
16
17 # Add edges to the graph
18 G.add_edges_from(unique_connections_df.values)
19
20 # Calculate degree for each node
21 degrees = G.degree()
22
23 # Get the 150 most friendly characters based on degree
24 top_characters = sorted(degrees, key=lambda x: x[1], reverse=True)[:10]
25
26 # Create a subgraph with the top characters
27 subgraph = G.subgraph([character for character, degree in top_characters])
28
29 # Get node types from nodes.csv
30 nodes = pd.read_csv('nodes.csv')
31
32 # Assuming 'node' column contains identifiers and 'type' column contains node types
33 node_types = dict(zip(nodes['node'], nodes['type']))
34
35 # Get unique node types including 'Unknown'
36 unique_node_types = list(set(node_types.values())) + ['Unknown']
37
38 # Get colors based on node types
39 color_palette = sns.color_palette("husl", n_colors=len(unique_node_types))
40 colors = [color_palette[unique_node_types.index(node_types.get(node, 'Unknown'))] for node in subgraph.nodes()]
41
42 # Draw the graph with node sizes based on degree and a single color for all edges
43 plt.figure(figsize=(12, 12))
44 nx.draw(subgraph, with_labels=True, node_size=[degrees[node] * 50 for node in subgraph.nodes()], node_color=colors, cmap=plt.cm.rainbow, edge_color='grey', alpha=0.7)
45 plt.title("Marvel Universe Social Network")
46 plt.show()
47
```

Marvel Universe Social Network



```
1 # Most 10 friendly not colored. only for easy demonstraion.
2
3 import pandas as pd
4 import networkx as nx
5 import matplotlib.pyplot as plt
6
7 # Load the hero network dataset
8 hero_network_df = pd.read_csv('hero-network.csv')
9
10 # Drop duplicate connections
11 unique_connections_df = hero_network_df[['hero1', 'hero2']].drop_duplicates()
12
13 # Create a graph object
14 G = nx.Graph()
15
16 # Add edges to the graph
17 G.add_edges_from(unique_connections_df.values)
18
19 # Calculate degree for each node
20 degrees = G.degree()
21
22 # Get the 150 most friendly characters based on degree
23 top_characters = sorted(degrees, key=lambda x: x[1], reverse=True)[:10]
24
25 # Create a subgraph with the top characters
26 subgraph = G.subgraph([character for character, degree in top_characters])
27
28 # Draw the graph with node sizes based on degree without any colors
29 plt.figure(figsize=(12, 12))
30 nx.draw(subgraph, with_labels=True, node_size=[degrees[node] * 50 for node in subgraph.nodes()], node_color='white', edge_color='grey', alpha=0.7)
31 plt.title("Marvel Universe Social Network")
32 plt.show()
```

Marvel Universe Social Network



```

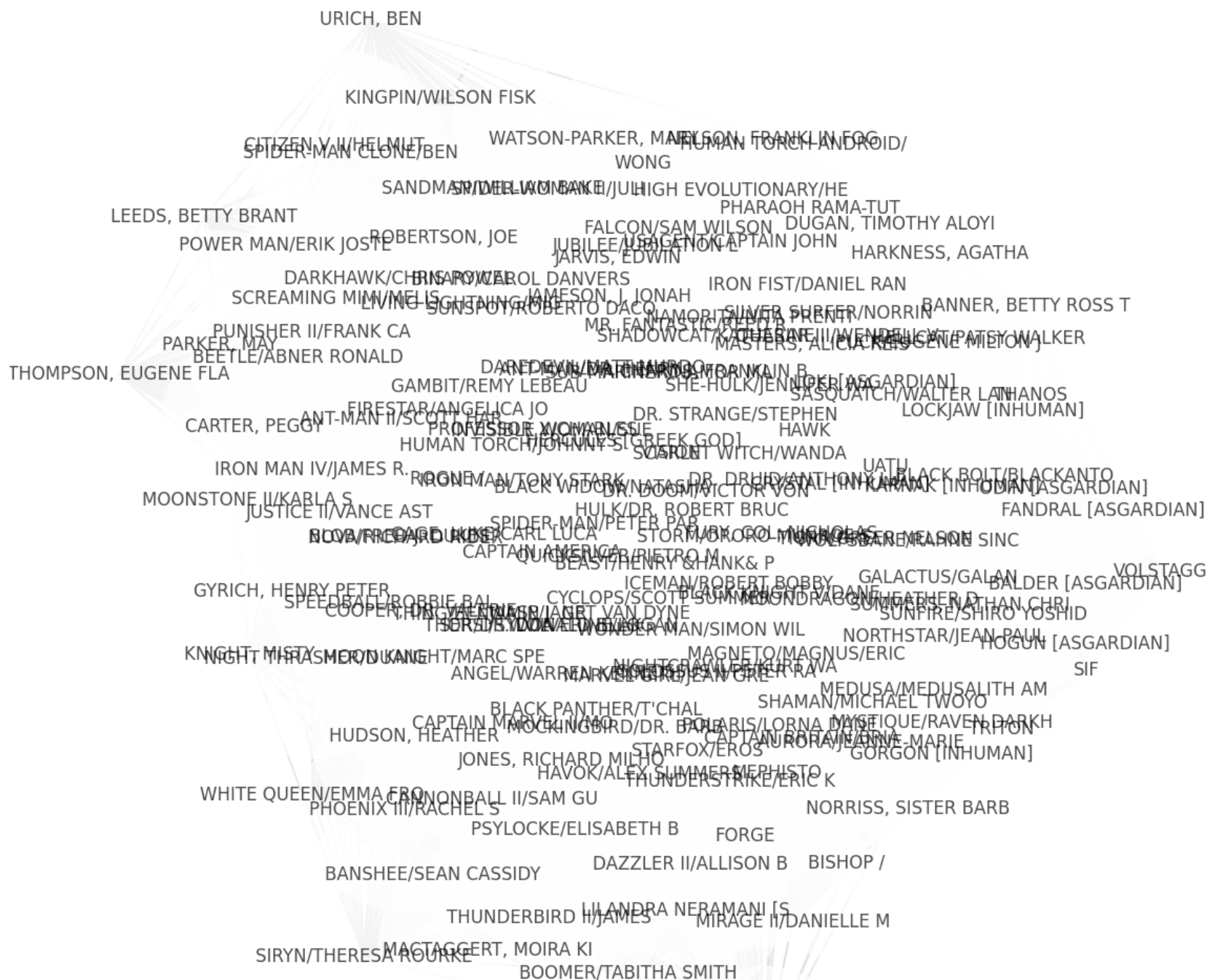
1 # Most 150 friendly colored.
2
3 import pandas as pd
4 import networkx as nx
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7
8 # Load the hero network dataset
9 hero_network_df = pd.read_csv('hero-network.csv')
10
11 # Drop duplicate connections
12 unique_connections_df = hero_network_df[['hero1', 'hero2']].drop_duplicates()
13
14 # Create a graph object
15 G = nx.Graph()
16
17 # Add edges to the graph
18 G.add_edges_from(unique_connections_df.values)
19
20 # Calculate degree for each node
21 degrees = G.degree()
22
23 # Get the 150 most friendly characters based on degree
24 top_characters = sorted(degrees, key=lambda x: x[1], reverse=True)[:150]
25
26 # Create a subgraph with the top characters
27 subgraph = G.subgraph([character for character, degree in top_characters])
28
29 # Get node types from nodes.csv
30 nodes = pd.read_csv('nodes.csv')
31
32 # Assuming 'node' column contains identifiers and 'type' column contains node types
33 node_types = dict(zip(nodes['node'], nodes['type']))
34
35 # Get unique node types including 'Unknown'
36 unique_node_types = list(set(node_types.values())) + ['Unknown']
37
38 # Get colors based on node types
39 color_palette = sns.color_palette("husl", n_colors=len(unique_node_types))
40 colors = [color_palette[unique_node_types.index(node_types.get(node, 'Unknown'))] for node in subgraph.nodes()]
41
42 # Draw the graph with node sizes based on degree and a single color for all edges
43 plt.figure(figsize=(12, 12))
44 nx.draw(subgraph, with_labels=True, node_size=[degrees[node] * 50 for node in subgraph.nodes()], node_color=colors, cmap=plt.cm.rainbow, edge_color='grey', alpha=0.7)
45 plt.title("Marvel Universe Social Network")
46 plt.show()
47

```



```
1 # Most 150 friendly not colored.
2
3 import pandas as pd
4 import networkx as nx
5 import matplotlib.pyplot as plt
6
7 # Load the hero network dataset
8 hero_network_df = pd.read_csv('hero-network.csv')
9
10 # Drop duplicate connections
11 unique_connections_df = hero_network_df[['hero1', 'hero2']].drop_duplicates()
12
13 # Create a graph object
14 G = nx.Graph()
15
16 # Add edges to the graph
17 G.add_edges_from(unique_connections_df.values)
18
19 # Calculate degree for each node
20 degrees = G.degree()
21
22 # Get the 150 most friendly characters based on degree
23 top_characters = sorted(degrees, key=lambda x: x[1], reverse=True)[:150]
24
25 # Create a subgraph with the top characters
26 subgraph = G.subgraph([character for character, degree in top_characters])
27
28 # Draw the graph with node sizes based on degree without any colors
29 plt.figure(figsize=(12, 12))
30 nx.draw(subgraph, with_labels=True, node_size=[degrees[node] * 50 for node in subgraph.nodes()], node_color='white', edge_color='grey', alpha=0.7)
31 plt.title("Marvel Universe Social Network")
32 plt.show()
```

Marvel Universe Social Network



Bonus: Visualize the above network using [Cytoscape](#) or [Gephi](#) or [D3](#) (10pt)

```
1 !pip install py4cytoscape
```

```
1 # I succeeded only when i downloaded Cytoscape manually on my pc.
```

```
1 _PY4CYTOSCAPE = 'git+https://github.com/cytoscape/py4cytoscape@0.0.11'
2 import requests
3 exec(requests.get("https://raw.githubusercontent.com/cytoscape/jupyter-bridge/master/client/p4c_init.py").text)
4 IPython.display.Javascript(_PY4CYTOSCAPE_BROWSER_CLIENT_JS) # Start browser client
```

```
1 import py4cytoscape as p4c
2 p4c.cytoscape_ping()
3 p4c.cytoscape_version_info()
```