

Part I. Text Analysis

A. Overview of Tools for Text Analysis

Some popular tools and libraries for text analysis include:

- **NLTK** (Natural Language Toolkit) – A Python library for processing textual data.
- **spaCy** – A fast, production-ready NLP library in Python.
- **TextBlob** – Simplified text processing for Python.
- **Gensim** – For topic modeling and similarity detection.
- **Stanford CoreNLP** – A suite of NLP tools written in Java.
- **Voyant Tools** – A web-based text reading and analysis environment.
- **AntConc** – A freeware corpus analysis toolkit.
- **LIWC** – A tool for linguistic inquiry and word count analysis.
- **Mallet** – For topic modeling and clustering.
- **RapidMiner/KNIME/Orange** – Platforms that offer text mining modules.

In this assignment, we will use **NLTK** for static and dynamic text analysis.

B. Static & Dynamic Text Analysis Example

1. Research Question / Application

Example Research Question:

“How does the frequency of common words vary across different novel sections?”

Application:

Understanding the distribution of keywords (e.g., common words or keywords of interest) across a text can help reveal shifts in themes, narrative focus, or stylistic changes.

2. Choose Method

- **Static Analysis:**
 - Calculate overall word frequency in the entire text.
 - Visualize the top N words using a bar plot.
- **Dynamic Analysis:**
 - Divide the text into segments (e.g., chapters or fixed-size blocks of sentences).
 - Track the frequency of a target word (or a set of words) across these segments.
 - Plot the frequency trends over segments to observe changes over time.

Expected Outcome:

- A **static summary** (frequency distribution of words) that shows the overall most common words in the text.
- A **dynamic visualization** (line plot) showing how the frequency of a specific word (or words) changes as you progress through the text.

3. Choose the Tool: NLTK

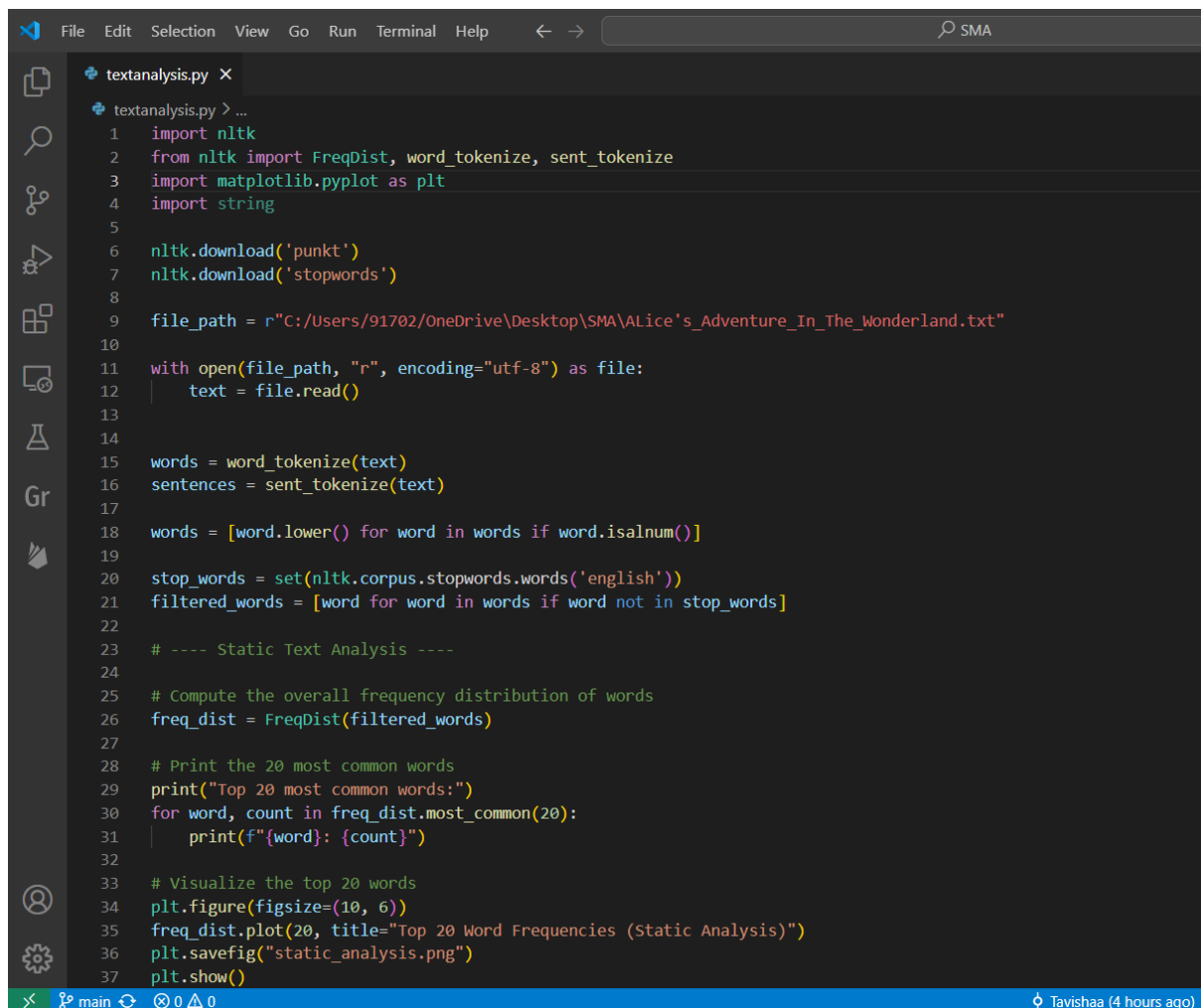
Why NLTK?

- **NLTK** is a widely used Python library for natural language processing.
- It offers powerful tokenization, frequency distribution, and text processing functions.
- It can be easily integrated with visualization libraries like Matplotlib to produce static and dynamic plots.

4. Do Analysis by Uploading Text Data

Text data link: <https://www.gutenberg.org/cache/epub/11/pg11.txt>

1) Static Text Analysis Code:

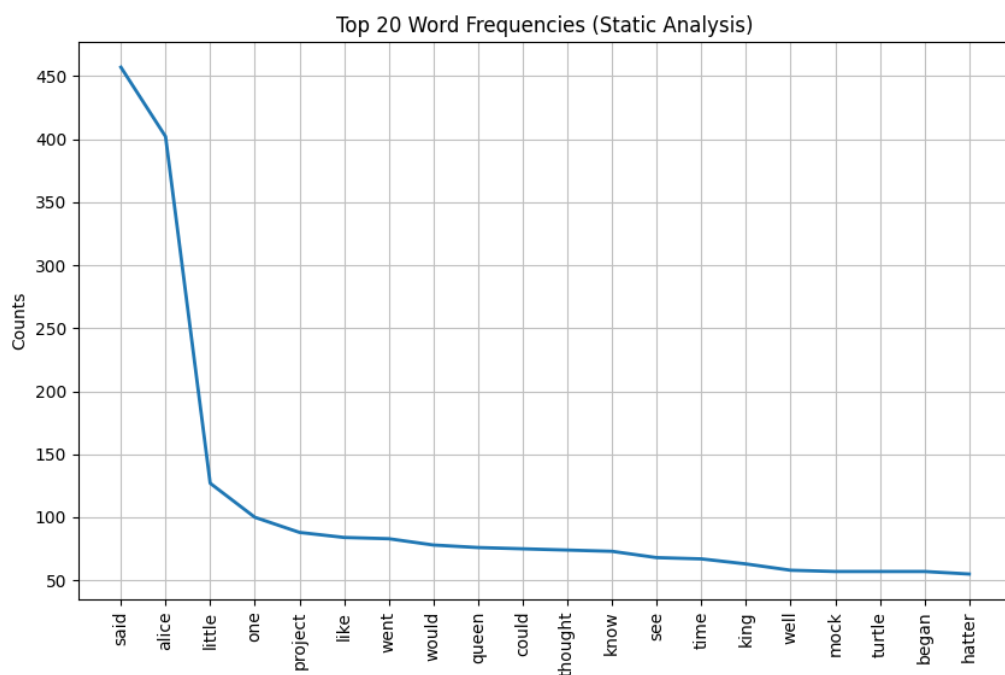


```
File Edit Selection View Go Run Terminal Help SMA
textanalysis.py X
textanalysis.py > ...
1 import nltk
2 from nltk import FreqDist, word_tokenize, sent_tokenize
3 import matplotlib.pyplot as plt
4 import string
5
6 nltk.download('punkt')
7 nltk.download('stopwords')
8
9 file_path = r"C:/Users/91702/OneDrive/Desktop/SMA/Alice's_Adventure_In_The_Wonderland.txt"
10
11 with open(file_path, "r", encoding="utf-8") as file:
12     text = file.read()
13
14
15 words = word_tokenize(text)
16 sentences = sent_tokenize(text)
17
18 words = [word.lower() for word in words if word.isalnum()]
19
20 stop_words = set(nltk.corpus.stopwords.words('english'))
21 filtered_words = [word for word in words if word not in stop_words]
22
23 # ---- Static Text Analysis ----
24
25 # Compute the overall frequency distribution of words
26 freq_dist = FreqDist(filtered_words)
27
28 # Print the 20 most common words
29 print("Top 20 most common words:")
30 for word, count in freq_dist.most_common(20):
31     print(f"{word}: {count}")
32
33 # Visualize the top 20 words
34 plt.figure(figsize=(10, 6))
35 freq_dist.plot(20, title="Top 20 Word Frequencies (Static Analysis)")
36 plt.savefig("static_analysis.png")
37 plt.show()
```

main 0 0 0 Tavishaa (4 hours ago)

Static Text Analysis Output (Most Common Words):

```
File Edit Selection View Go Run Terminal Help
PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS
PS C:\Users\91702\OneDrive\Desktop\SMA> python textanalysis.py
[nltk_data] Downloading package punkt to
[nltk_data]   C:\Users\91702\AppData\Roaming\nltk_data...
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package stopwords to
[nltk_data]   C:\Users\91702\AppData\Roaming\nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
Top 20 most common words:
said: 457
alice: 402
little: 127
one: 100
project: 88
like: 84
went: 83
would: 78
queen: 76
could: 75
thought: 74
know: 73
see: 68
time: 67
king: 63
well: 58
mock: 57
turtle: 57
began: 57
hatter: 55
```

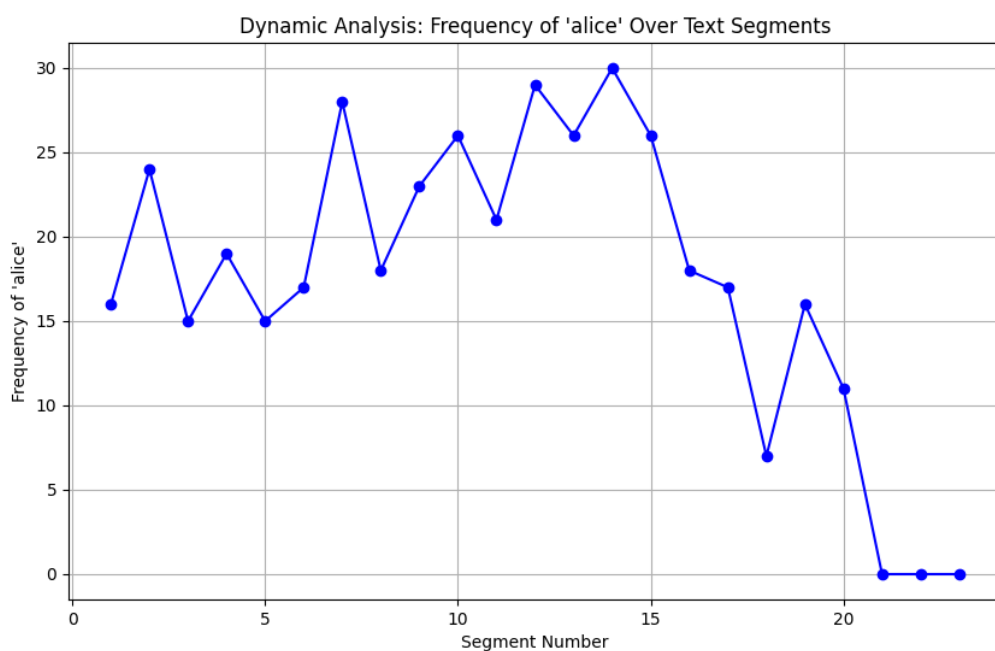


2) Dynamic Text Analysis Code:

```
File Edit Selection View Go Run Terminal Help SMA
textanalysis.py X
textanalysis.py > ...

39 # ---- Dynamic Text Analysis ----
40
41 # Analyze how a target word changes across the text
42 segment_size = 50
43 num_sentences = len(sentences)
44 segments = [sentences[i:i + segment_size] for i in range(0, num_sentences, segment_size)]
45
46 # Choose a target word to track
47 target_word = "alice"
48 dynamic_freq = []
49
50 # Calculate the frequency of the target word in each segment
51 for idx, segment in enumerate(segments, 1):
52     segment_text = " ".join(segment)
53     segment_words = word_tokenize(segment_text)
54     segment_words = [word.lower() for word in segment_words if word.isalnum()] # Clean words
55     count = segment_words.count(target_word)
56     dynamic_freq.append(count)
57     print(f"Segment {idx}: Frequency of '{target_word}' = {count}")
58
59 # Plot the frequency of the target word over segments
60 plt.figure(figsize=(10, 6))
61 plt.plot(range(1, len(dynamic_freq) + 1), dynamic_freq, marker='o', linestyle='-', color='b')
62 plt.xlabel("Segment Number")
63 plt.ylabel(f"Frequency of '{target_word}'")
64 plt.title(f"Dynamic Analysis: Frequency of '{target_word}' Over Text Segments")
65 plt.grid(True)
66 plt.savefig("dynamic_analysis.png")
67 plt.show()
68
```

Dynamic Text Analysis Output:



Part II. Hyperlink Analysis

A. Overview of Tools for Hyperlink Analysis

Some popular tools and software for hyperlink and network analysis include:

- **Gephi** – An open-source network visualization and analysis tool.
- **NodeXL** – A network analysis tool that works as an Excel add-in.
- **Cytoscape** – A software platform for visualizing complex networks.
- **Pajek** – A program for large network analysis.
- **UCINET** – Software for social network analysis.
- **SocNetV** – A tool for social network visualization.
- **Netlytic** – For analyzing social media networks.
- **Graph-tool** – A Python library for network analysis.
- **Linkurious** and **Maltego** – For investigative and forensic link analysis.

B. Static & Dynamic Hyperlink Analysis Example

1. Purpose of the Analysis

Hyperlink analysis is used to examine the structure of web links to gain insights into:

- Identifying **influential webpages** within Wikipedia
- Understanding the **interlinking structure** of Wikipedia articles
- Mapping **the most referenced Wikipedia pages**
- Analyzing **how topics are connected** through hyperlinks

2. Defining the Network

A hyperlink network consists of **nodes (Wikipedia pages)** and **edges (hyperlinks between them)**. It can be represented as:

- **Directed Graph:** Links have direction (e.g., Page A → Page B means Page A links to Page B)
- **Weighted Graph:** Some links have more importance than others based on frequency
- **Community Graph:** Groups of pages closely linked together

3. Tools for Hyperlink Analysis

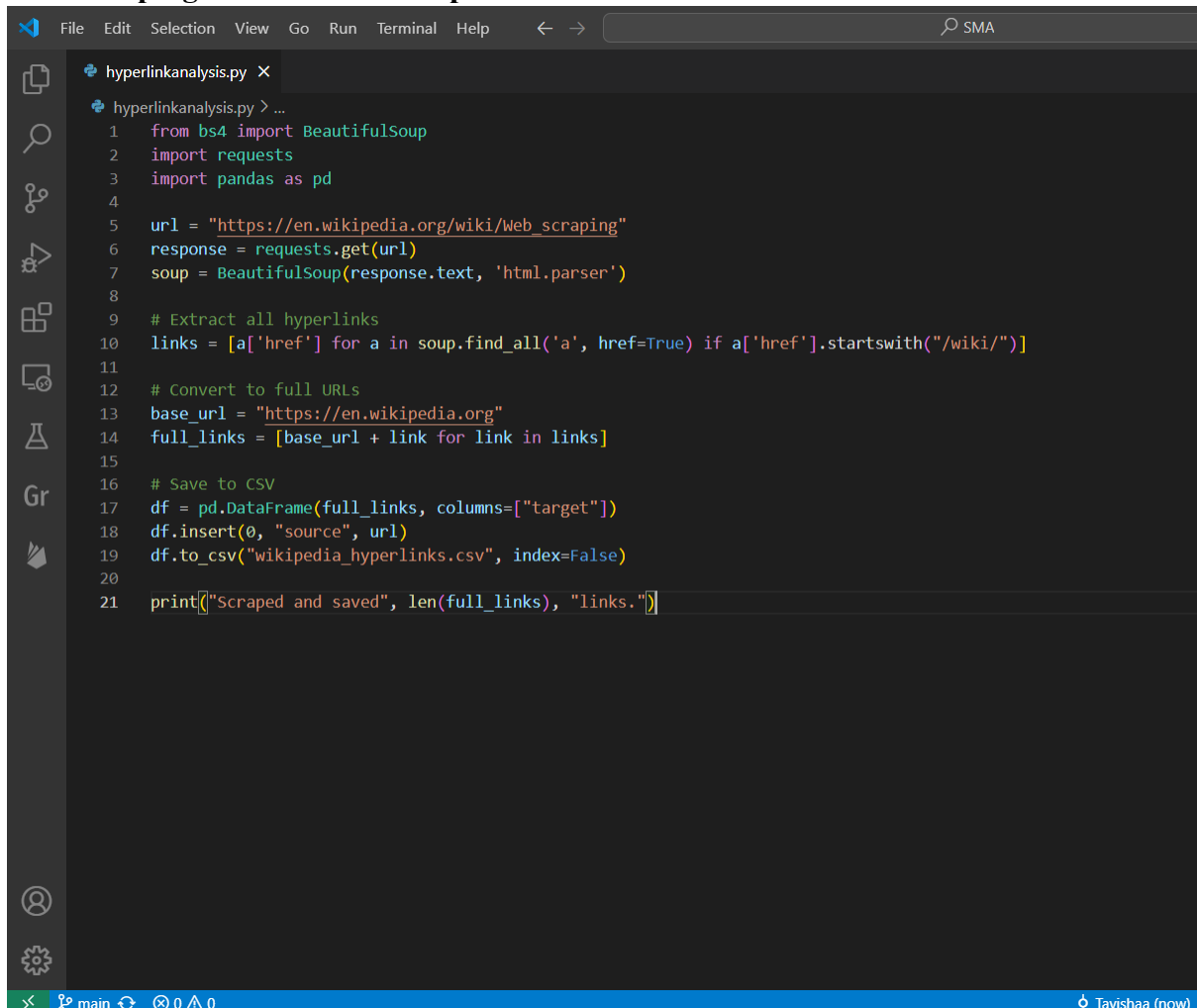
For this assignment, we will use:

- **BeautifulSoup (Python):** To scrape hyperlinks from Wikipedia
- **NetworkX (Python):** To build and analyze the hyperlink network
- **Matplotlib:** To visualize the network structure

4. Scraping Wikipedia Hyperlinks

We will scrape hyperlinks from the Wikipedia page on **Web Scraping** (https://en.wikipedia.org/wiki/Web_scraping).

Web Scraping with BeautifulSoup:



The screenshot shows a Visual Studio Code editor window with a file named `hyperlinkanalysis.py`. The script uses `BeautifulSoup` to scrape hyperlinks from the Wikipedia page on Web Scraping. The code is as follows:

```
1 from bs4 import BeautifulSoup
2 import requests
3 import pandas as pd
4
5 url = "https://en.wikipedia.org/wiki/Web_scraping"
6 response = requests.get(url)
7 soup = BeautifulSoup(response.text, 'html.parser')
8
9 # Extract all hyperlinks
10 links = [a['href'] for a in soup.find_all('a', href=True) if a['href'].startswith("/wiki/")]
11
12 # Convert to full URLs
13 base_url = "https://en.wikipedia.org"
14 full_links = [base_url + link for link in links]
15
16 # Save to CSV
17 df = pd.DataFrame(full_links, columns=["target"])
18 df.insert(0, "source", url)
19 df.to_csv("wikipedia_hyperlinks.csv", index=False)
20
21 print("Scraped and saved", len(full_links), "links.")
```

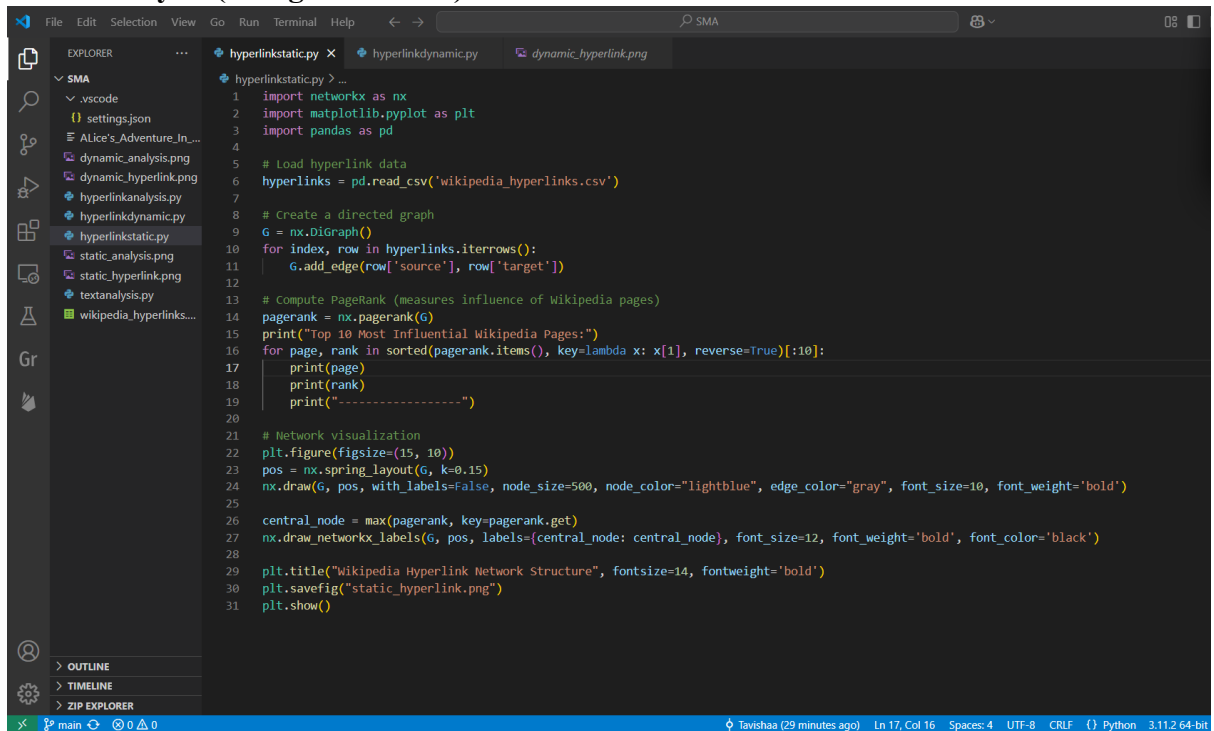
The editor interface includes a sidebar with icons for Explorer, Search, Source Control, Run and Debug, Extensions, and a bottom status bar showing the current branch as `main` and the user as `Tavishaa (now)`.

Link for the csv file created:

https://github.com/Tavishaa/SMA-analysis/blob/main/wikipedia_hyperlinks.csv

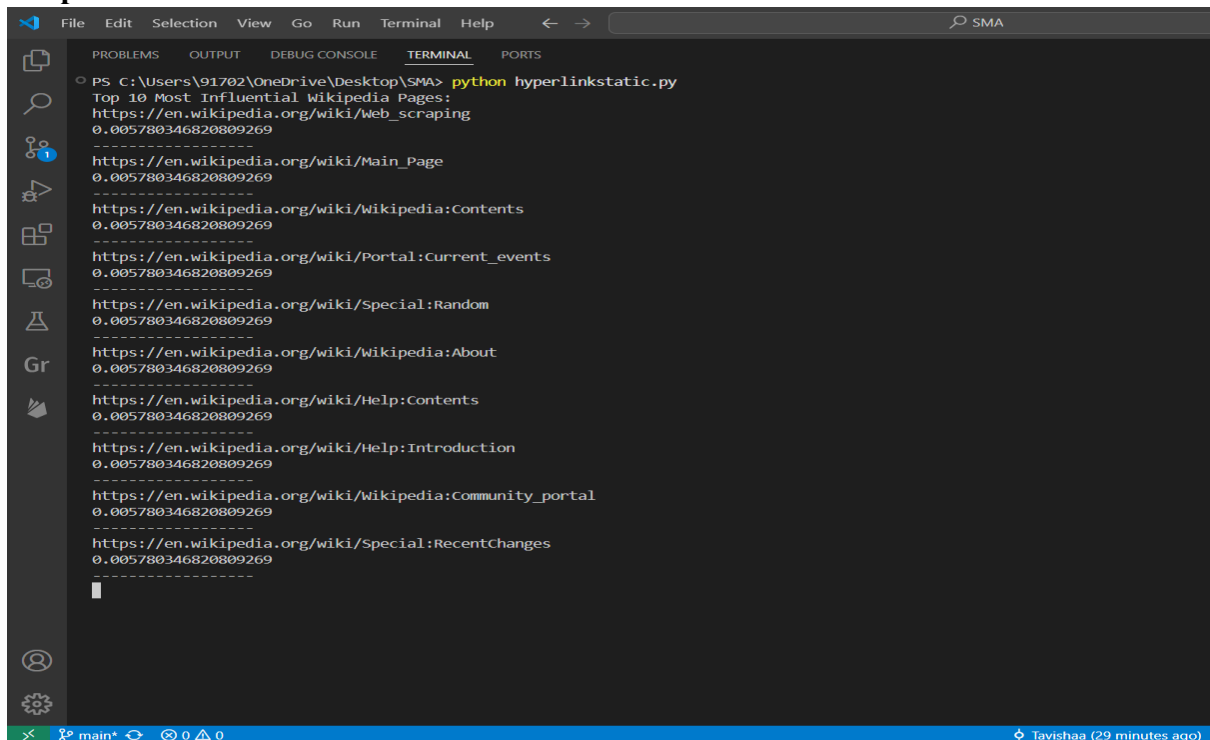
5. Performing Static and Dynamic Analysis

Static Analysis (Using NetworkX)

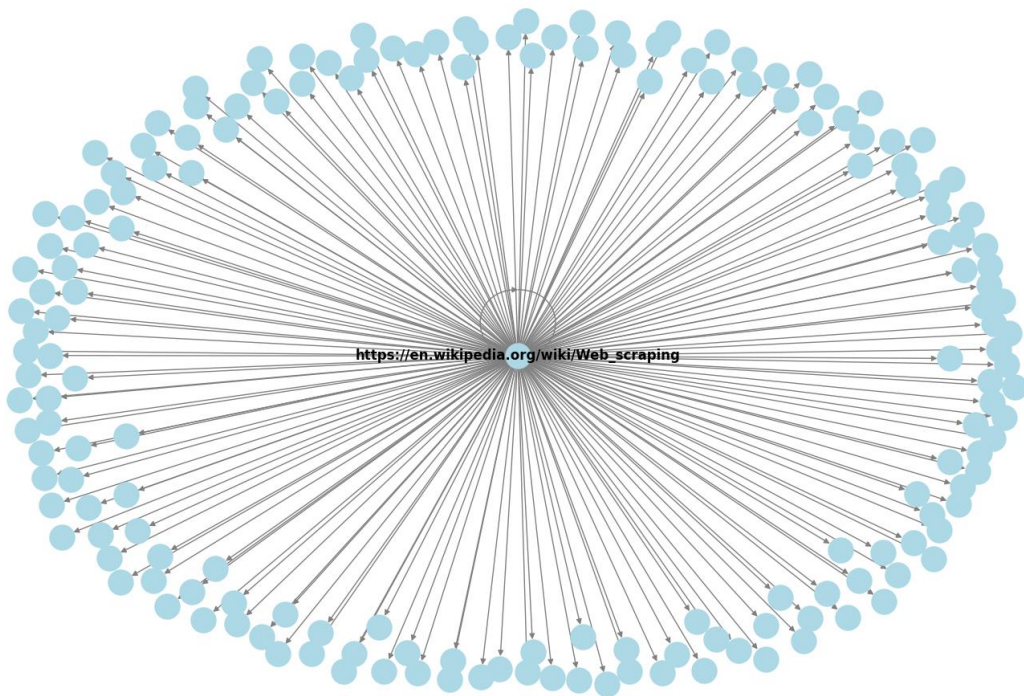


```
1 import networkx as nx
2 import matplotlib.pyplot as plt
3 import pandas as pd
4
5 # Load hyperlink data
6 hyperlinks = pd.read_csv('wikipedia_hyperlinks.csv')
7
8 # Create a directed graph
9 G = nx.DiGraph()
10 for index, row in hyperlinks.iterrows():
11     G.add_edge(row['source'], row['target'])
12
13 # Compute PageRank (measures influence of Wikipedia pages)
14 pagerank = nx.pagerank(G)
15 print("Top 10 Most Influential Wikipedia Pages:")
16 for page, rank in sorted(pagerank.items(), key=lambda x: x[1], reverse=True)[:10]:
17     print(page)
18     print(rank)
19     print("-----")
20
21 # Network visualization
22 plt.figure(figsize=(15, 10))
23 pos = nx.spring_layout(G, k=0.15)
24 nx.draw(G, pos, with_labels=False, node_size=500, node_color="lightblue", edge_color="gray", font_size=10, font_weight='bold')
25
26 central_node = max(pagerank, key=pagerank.get)
27 nx.draw_networkx_labels(G, pos, labels={central_node: central_node}, font_size=12, font_weight='bold', font_color='black')
28
29 plt.title("Wikipedia Hyperlink Network Structure", fontsize=14, fontweight='bold')
30 plt.savefig("static_hyperlink.png")
31 plt.show()
```

Output:



```
PS C:\Users\91702\OneDrive\Desktop\SMA> python hyperlinkstatic.py
Top 10 Most Influential Wikipedia Pages:
https://en.wikipedia.org/wiki/Web_scraping
0.005780346820809269
-----
https://en.wikipedia.org/wiki/Main_Page
0.005780346820809269
-----
https://en.wikipedia.org/wiki/Wikipedia:Contents
0.005780346820809269
-----
https://en.wikipedia.org/wiki/Portal:Current_events
0.005780346820809269
-----
https://en.wikipedia.org/wiki/Special:Random
0.005780346820809269
-----
https://en.wikipedia.org/wiki/Wikipedia:About
0.005780346820809269
-----
https://en.wikipedia.org/wiki/Help:Contents
0.005780346820809269
-----
https://en.wikipedia.org/wiki/Help:Introduction
0.005780346820809269
-----
https://en.wikipedia.org/wiki/Wikipedia:Community_portal
0.005780346820809269
-----
https://en.wikipedia.org/wiki/Special:RecentChanges
0.005780346820809269
-----
```



Dynamic Analysis (Tracking Changes Over Time)

```

hyperlinkdynamic.py X
hyperlinkdynamic.py > ...
1 import networkx as nx
2 import matplotlib.pyplot as plt
3
4 # Function to update network over time
5 def update_network(graph, new_links):
6     for src, tgt in new_links:
7         graph.add_edge(src, tgt)
8     return graph
9
10 # Ensure 'G' is defined before adding new links
11 if 'G' not in globals():
12     G = nx.DiGraph()
13
14 # Simulating new links being added dynamically
15 new_links = [("https://en.wikipedia.org/wiki/Web_scraping", "https://en.wikipedia.org/wiki/Data_scraping"),
16             ("https://en.wikipedia.org/wiki/Web_scraping", "https://en.wikipedia.org/wiki/Web_crawler")]
17 G = update_network(G, new_links)
18
19 # Recompute PageRank dynamically
20 pagerank_updated = nx.pagerank(G)
21 print("Updated Top 10 Most Influential Wikipedia Pages:")
22 for page, rank in sorted(pagerank_updated.items(), key=lambda x: x[1], reverse=True)[:10]:
23     print(page, "\n", rank, "\n-----")
24
25 # Improved Graph Visualization
26 plt.figure(figsize=(15, 10))
27 pos = nx.spring_layout(G, k=0.3) # Adjust k for better spacing
28 nx.draw(G, pos, with_labels=False, node_size=500, node_color="lightcoral", edge_color="gray", alpha=0.6)
29
30 # Label only the central node in the updated graph
31 central_node_updated = max(pagerank_updated, key=pagerank_updated.get)
32 nx.draw_networkx_labels(G, pos, labels={central_node_updated: central_node_updated},
33                        font_size=10, font_weight='bold', font_color='red',
34                        bbox=dict(facecolor='white', edgecolor='black', boxstyle='round,pad=0.3'))
35
36 plt.title("Updated Wikipedia Hyperlink Network Structure", fontsize=14, fontweight='bold')
37 plt.savefig("dynamic_hyperlink.png", bbox_inches="tight") # Save graph properly

```


Output:

```
File Edit Selection View Go Run Terminal Help SMA
PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS
PS C:\Users\91702\OneDrive\Desktop\SMA> python hyperlinkdynamic.py
Updated Top 10 Most Influential Wikipedia Pages:
https://en.wikipedia.org/wiki/Data_scraping
0.37012974744707666
-----
https://en.wikipedia.org/wiki/Web_crawler
0.37012974744707666
-----
https://en.wikipedia.org/wiki/Web_scraping
0.25974050510584634
-----
PS C:\Users\91702\OneDrive\Desktop\SMA>
```

Updated Wikipedia Hyperlink Network Structure

