

Krishna Khandelwal

Assignmnet-2

Assignment 2: NLP Analysis of News Articles Using the NYTimes API

Background

Natural Language Processing (NLP) plays a crucial role in extracting insights from large volumes of textual data. News articles contain valuable information about public sentiment, social trends, and emerging topics. Businesses, policymakers, and researchers use NLP to analyze trends in various domains, such as finance, politics, technology, and public health.

Industries utilize news analytics for:

- **Market Intelligence** – Companies monitor news for insights into competitors, industry trends, and economic conditions.
- **Public Sentiment & Brand Monitoring** – Organizations track public opinion about key figures (e.g., CEOs) or events.
- **Crisis Management** – Analyzing real-time news helps businesses manage PR crises and understand public perception.
- **Regulatory & Compliance Monitoring** – Financial and legal industries track regulatory changes and enforcement actions.
- **Risk Analysis** – Financial institutions and investors use news-based sentiment to predict market fluctuations.

This assignment will guide you through **collecting, processing, and analyzing news articles** to extract meaningful insights.

✓ Step 1: Build a Dataset Using the NYTimes API

```
1 import requests as req
2 import time
3 import pandas as pd
4 from tqdm import tqdm
5
6 # User inputs
7 TOPIC = input("Enter the topic: ").strip()
8 START_DATE = input("Enter start date (YYYYMMDD): ").strip()
9 END_DATE = input("Enter end date (YYYYMMDD): ").strip()
```

```
10 TOTAL_ARTICLES = int(input("Enter the total number of articles to fetch: ").strip())
11
12 # API Key (Replace with your valid key)
13 API_KEY = "TFI8GCddEI7CWkkZJuAiCz50NeVbLvev"
14
15 # Base URL for NY Times API
16 BASE_URL = "https://api.nytimes.com/svc/search/v2/articlesearch.json"
17
18 # Storage for articles
19 articles = []
20 pages = TOTAL_ARTICLES // 10 + (1 if TOTAL_ARTICLES % 10 > 0 else 0)
21
22 # Fetching articles
23 for page in tqdm(range(pages), desc="Fetching Articles"):
24     url = (
25         f"{BASE_URL}?q={TOPIC}&begin_date={START_DATE}&end_date={END_DATE}"
26         f"&api-key={API_KEY}&page={page}"
27     )
28
29     response = req.get(url).json()
30
31     if "response" in response and "docs" in response["response"]:
32         docs = response["response"]["docs"]
33
34         for doc in docs:
35             articles.append({
36                 "Title": doc["headline"]["main"] if "headline" in doc else None,
37                 "Publication Date": doc.get("pub_date", None),
38                 "Author": doc.get("byline", {}).get("original", None),
39                 "Location": doc.get("section_name", None),
40                 "Summary": doc.get("abstract", None),
41                 "Full Article Content": doc.get("lead_paragraph", None),
42                 "URL": doc.get("web_url", None),
43             })
44
45         if len(articles) >= TOTAL_ARTICLES:
46             break
47
48     time.sleep(6)
49
50 # Convert to DataFrame
51 df = pd.DataFrame(articles)
52
53 # Save to CSV
54 csv_filename = "nytimes_articles.csv"
55 df.to_csv(csv_filename, index=False, encoding="utf-8")
56
57 print(f"\nDataset saved as {csv_filename}")
58
```

```
→ Enter the topic: Chatgpt
Enter start date (YYYYMMDD): 20150101
Enter end date (YYYYMMDD): 20250225
Enter the total number of articles to fetch: 100
Fetching Articles: 90%|██████████| 9/10 [01:09<00:07, 7.67s/it]
Dataset saved as nytimes_articles.csv
```

Collected news articles from the New York Times Article Search API based on chatgpt topic and date range provided by the user. It will retrieves key details like the title, author, and summary, then saves the results into a CSV file for further analysis.

```
1 len(articles)
```

```
→ 100
```

Checked for length of the articles which i tried to fetch

All necessary libraries

```
1 import pandas as pd
2 import plotly.express as px
3 import plotly.graph_objects as go
4
5 # Install necessary libraries
6 !pip install spacy
7 !python -m spacy download en_core_web_sm
8
9 # Import required libraries
10 import spacy
11 import re
12 import string
13 from nltk.corpus import stopwords
14 import nltk
15
16 # Install VADER for sentiment analysis
17 !pip install vaderSentiment
18
19 # Import necessary libraries
20 import pandas as pd
21 import matplotlib.pyplot as plt
22 from vaderSentiment.vaderSentiment import SentimentIntensityAnalyzer
23
24 # Install wordcloud if not already installed
25 !pip install wordcloud
26
27 # Import necessary libraries
```

```
28 import pandas as pd
29 import matplotlib.pyplot as plt
30 from wordcloud import WordCloud
31
32 # Install necessary libraries if not already installed
33 !pip install gensim
34
35 # Import necessary libraries
36 import pandas as pd
37 import gensim
38 from gensim import corpora
39 from gensim.models import LdaModel
40 from nltk.tokenize import word_tokenize
41 import nltk
42 import matplotlib.pyplot as plt
43 import nltk
44 nltk.download('punkt')
45
46
47 nltk.download('punkt_tab')
48
49
50
```

▼ Step 2: Save Data as a CSV File

```
1 import pandas as pd
2
3 # Load the saved dataset
4 df = pd.read_csv("nytimes_articles.csv")
5
6 # Display the first few rows
7 print(df.head(10))
8
9 # Check for missing values
10 print(df.isnull().sum())
11
```

→ Title \

0	She Fell in Love With ChatGPT. Like, Actual Lo...
1	ChatGPT Is Restored After Hourslong Outage
2	She Is in Love With ChatGPT
3	How Helpful Is Operator, OpenAI's New A.I. Agent?
4	They Invested Billions. Then the A.I. Script G...
5	First Impressions of DeepSeek's A.I. Chatbot
6	OpenAI Unveils A.I. Agent That Can Use Website...
7	OpenAI Courts Trump With Vision for 'A.I. in A...
8	At the Intersection of A.I. and Spirituality
9	The 2024 Good Tech Awards

```
Publication Date \
0 2025-02-25T11:00:10+0000
1 2025-01-23T13:02:14+0000
2 2025-01-15T14:00:09+0000
3 2025-02-01T10:01:31+0000
4 2025-01-29T10:03:18+0000
5 2025-01-27T23:33:09+0000
6 2025-01-23T18:00:06+0000
7 2025-01-13T11:00:10+0000
8 2025-01-03T10:00:23+0000
9 2024-12-30T17:00:02+0000
```

	Author	Location	\
0	By Natalie Kitroeff, Kashmir Hill, Nina Feldma...	Podcasts	
1	By Victor Mather	Business Day	
2	By Kashmir Hill	Technology	
3	By Kevin Roose	Technology	
4	By Erin Griffith	Technology	
5	By Eli Tan	Technology	
6	By Cade Metz	Technology	
7	By Cade Metz and Cecilia Kang	Technology	
8	By Eli Tan	Technology	
9	By Kevin Roose	Technology	

```
Summary \
0 The story of a woman who is dating an A.I. cha...
1 Thousands of users were unable to access OpenA...
2 A 28-year-old woman with a busy social life sp...
3 Operator, a new computer-using tool from OpenA...
4 Venture capitalists plowed money into A.I. sta...
5 The chatbot from China appears to perform a nu...
6 The new tool, called Operator, can shop for gr...
7 The maker of ChatGPT hopes to spur investment ...
8 Modern religious leaders are experimenting wit...
9 In a year of continued A.I. progress, “founder...
```

```
Full Article Content \
0 Warning: This episode discusses sexual themes.
1 ChatGPT was disrupted for several hours on Thu...
2 Ayrin's love affair with her A.I. boyfriend st...
3 In the past week, OpenAI's Operator has done t...
4 Jordan Jacobs, an investor at the venture capi...
5 A new chatbot created by the Chinese company D...
6 Two years ago, OpenAI launched the chatbot cra...
7 In December, Sam Altman, OpenAI's chief execut...
```

After collecting the articles, I saved them into a CSV file called nytimes_articles.csv for easy access and further analysis. In this step, I loaded the dataset back into a pandas DataFrame, previewed the first few rows to check the data structure, and verified if there were any missing values. This helped ensure the dataset was clean and ready for the next steps.

▼ Step 3: Load the Saved Dataset

```
1 # Performed basic inspection
2 df_info = df.info() # Get dataset info
3 missing_values = df.isnull().sum() # Count missing values per column
4 duplicate_rows = df.duplicated().sum() # Count duplicate rows
5
6

→ <class 'pandas.core.frame.DataFrame'>
RangeIndex: 100 entries, 0 to 99
Data columns (total 7 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Title            100 non-null    object  
 1   Publication Date 100 non-null    object  
 2   Author           99 non-null    object  
 3   Location          99 non-null    object  
 4   Summary           100 non-null    object  
 5   Full Article Content 100 non-null    object  
 6   URL              100 non-null    object  
dtypes: object(7)
memory usage: 5.6+ KB
```

```
1 # Display first few rows of the dataset
2 print(df.head())
3
4 # Display missing values
5 print("\nMissing values per column:")
6 print(df.isnull().sum())
7
8 # Display duplicate row count
9 print("\nDuplicate rows in dataset:", df.duplicated().sum())
10
11 # Display dataset information
12 print("\nDataset Info:")
13 df.info()
14
```

```
→                                     Title \
0  She Fell in Love With ChatGPT. Like, Actual Lo...
1      ChatGPT Is Restored After Hourslong Outage
2                      She Is in Love With ChatGPT
3  How Helpful Is Operator, OpenAI's New A.I. Agent?
4  They Invested Billions. Then the A.I. Script G...
```

```
                                     Publication Date \
0  2025-02-25T11:00:10+0000
1  2025-01-23T13:02:14+0000
2  2025-01-15T14:00:09+0000
3  2025-02-01T10:01:31+0000
```

4 2025-01-29T10:03:18+0000

	Author	Location	\
0	By Natalie Kitroeff, Kashmir Hill, Nina Feldma...	Podcasts	
1	By Victor Mather	Business Day	
2	By Kashmir Hill	Technology	
3	By Kevin Roose	Technology	
4	By Erin Griffith	Technology	

	Summary	\
0	The story of a woman who is dating an A.I. cha...	
1	Thousands of users were unable to access OpenA...	
2	A 28-year-old woman with a busy social life sp...	
3	Operator, a new computer-using tool from OpenA...	
4	Venture capitalists plowed money into A.I. sta...	

	Full Article Content	\
0	Warning: This episode discusses sexual themes.	
1	ChatGPT was disrupted for several hours on Thu...	
2	Ayrin's love affair with her A.I. boyfriend st...	
3	In the past week, OpenAI's Operator has done t...	
4	Jordan Jacobs, an investor at the venture capi...	

	URL
0	https://www.nytimes.com/2025/02/25/podcasts/th...
1	https://www.nytimes.com/2025/01/23/business/ch...
2	https://www.nytimes.com/2025/01/15/technology/...
3	https://www.nytimes.com/2025/02/01/technology/...
4	https://www.nytimes.com/2025/01/29/technology/...

Missing values per column:

Title	0
Publication Date	0
Author	1
Location	1
Summary	0
Full Article Content	0
URL	0

dtype: int64

Duplicate rows in dataset: 0

Dataset Info:

<class 'pandas.core.frame.DataFrame'>

RangeIndex: 100 entries, 0 to 99

Data columns (total 7 columns):

I inspected the dataset by displaying the first few rows and checking for missing values and duplicates. Then, I used df.info() to review the dataset's structure, ensuring it was clean and ready for analysis.

▼ Step 4: Exploratory Data Analysis (EDA)

Feature Engineering

```

1 # Ensure Full Article Content column is not empty
2 df["Full Article Content"] = df["Full Article Content"].fillna("")
3
4 # Part 1: Word Count per Article
5 df["Word Count"] = df["Full Article Content"].apply(lambda x: len(x.split()))
6
7 # Part 2: Sentence Count per Article
8 df["Sentence Count"] = df["Full Article Content"].apply(lambda x: len(x.split('.')))
9
10 # Part 3: Character Count per Article
11 df["Character Count"] = df["Full Article Content"].apply(lambda x: len(x))
12
13 # Part 4: Average Word Length per Article
14 df["Avg Word Length"] = df["Full Article Content"].apply(
15     lambda x: sum(len(word) for word in x.split()) / len(x.split()) if len(x.split()) >
16 )
17

```

In this step, I engineered new features to better understand the articles. I calculated the word count, sentence count, character count, and average word length for each article. These features helped summarize the length and complexity of the articles, providing useful insights for further analysis.

Some Visualizations

Article frequency over time

```

1 import plotly.express as px
2
3 # Convert 'Publication Date' to datetime (if not already converted)
4 df["Publication Date"] = pd.to_datetime(df["Publication Date"], errors="coerce")
5
6 # Extract year-month for time-based analysis
7 df["Year-Month"] = df["Publication Date"].dt.to_period("M").astype(str)
8
9 # Create a DataFrame for plotting
10 article_counts = df["Year-Month"].value_counts().sort_index().reset_index()
11 article_counts.columns = ["Year-Month", "Number of Articles"]
12
13 # Create bar chart
14 fig = px.bar(
15     article_counts,
16     x="Year-Month",
17     y="Number of Articles",

```

```
18     title="Article Frequency Over Time",
19     labels={"Year-Month": "Year-Month", "Number of Articles": "Number of Articles"},
20     color_discrete_sequence=["orange"]
21 )
22
23 fig.update_layout(
24     xaxis_tickangle=-45,
25     bargap=0.2,
26     plot_bgcolor='white'
27 )
28
29 fig.show()
30
```



I analyzed how many articles were published each month by grouping them by year and month. Then, I visualized this data using a bar chart in Plotly, which shows the frequency of articles over time. This helps identify trends in publication volume and highlights periods of high or low activity.

Distribution of article lengths

```
1 # Created histogram
2 fig = px.histogram(
3     df,
4     x="Word Count",
5     nbins=30,
6     title="Distribution of Article Lengths",
7     labels={"Word Count": "Word Count", "count": "Number of Articles"},
8     color_discrete_sequence=["orange"]
9 )
10
11 fig.update_layout(
12     bargap=0.1,
13     plot_bgcolor='white'
14 )
15
16 fig.show()
17
```



I created a histogram to visualize the distribution of article lengths based on word count. This helped me understand how long most articles are and identify whether the dataset is balanced between shorter and longer pieces.

Authors with the most articles

```
1 import plotly.graph_objects as go
2
3 # Get the top 10 authors by article count
4 top_authors = df["Author"].value_counts().head(10).reset_index()
5 top_authors.columns = ["Author", "Number of Articles"]
6
7 # Create vertical bar chart
8 fig = go.Figure(go.Bar(
9     x=top_authors["Author"],
10    y=top_authors["Number of Articles"],
11    text=top_authors["Number of Articles"],
12    textposition='outside',
13    marker=dict(
14        color='rgba(26, 118, 255, 0.7)',
15        line=dict(color='rgba(26, 118, 255, 1.0)', width=1.5)
16    )
17 ))
18
19 # Update layout
20 fig.update_layout(
21     title="Top Authors with Most Articles",
22     xaxis_title="Author",
23     yaxis_title="Number of Articles",
24     xaxis_tickangle=-15, # Slight tilt, easier to read
25     xaxis=dict(
26         tickfont=dict(size=12),
27         tickmode='array',
28         tickvals=list(range(len(top_authors))),
29         ticktext=[name[:30] + '...' if len(name) > 30 else name for name in top_authors[
30 ]),
31     yaxis=dict(tickfont=dict(size=12)),
32     title_font=dict(size=20),
33     plot_bgcolor='white',
34     width=900,
35     height=600
36 )
37
38 # Show the plot
39 fig.show()
40
```



I created a vertical bar chart to display the top 10 authors who published the most articles in the dataset. This visualization highlights the key contributors and gives a quick overview of who wrote the most on the selected topics. I also adjusted the layout and label formatting to make the chart more readable.

▼ Step 5: NLP Preprocessing Pipeline

```
1 # Install necessary libraries
2 !pip install spacy
3 !python -m spacy download en_core_web_sm
4
5 # Import required libraries
6 import spacy
```

```
7 import re
8 import string
9 from nltk.corpus import stopwords
10 import nltk
11
12 # Download stopwords
13 nltk.download('stopwords')
14
15 # Load SpaCy model
16 nlp = spacy.load("en_core_web_sm")
17
18
```

→ Requirement already satisfied: spacy in /usr/local/lib/python3.11/dist-packages (3.7.1)
Requirement already satisfied: spacy-legacy<3.1.0,>=3.0.11 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: spacy-loggers<2.0.0,>=1.0.0 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: murmurhash<1.1.0,>=0.28.0 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: cymem<2.1.0,>=2.0.2 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: preshed<3.1.0,>=3.0.2 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: thinc<8.3.0,>=8.2.2 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: wasabi<1.2.0,>=0.9.1 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: srsly<3.0.0,>=2.4.3 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: catalogue<2.1.0,>=2.0.6 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: weasel<0.5.0,>=0.1.0 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: typer<1.0.0,>=0.3.0 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: tqdm<5.0.0,>=4.38.0 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: requests<3.0.0,>=2.13.0 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: pydantic!=1.8,!<1.8.1,<3.0.0,>=1.7.4 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: jinja2 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: setuptools in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: langcodes<4.0.0,>=3.2.0 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: numpy>=1.19.0 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: language-data>=1.2 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: annotated-types>=0.6.0 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: pydantic-core==2.27.2 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: typing-extensions>=4.12.2 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: blis<0.8.0,>=0.7.8 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: confection<1.0.0,>=0.0.1 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: click>=8.0.0 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: shellingham>=1.3.0 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: rich>=10.11.0 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: cloudpathlib<1.0.0,>=0.7.0 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: smart-open<8.0.0,>=5.2.1 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: marisa-trie>=1.1.0 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: markdown-it-py>=2.2.0 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: pygments<3.0.0,>=2.13.0 in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: wrapt in /usr/local/lib/python3.11/dist-packages (from spacy)
Requirement already satisfied: mdurl~>=0.1 in /usr/local/lib/python3.11/dist-packages (from spacy)
Collecting en-core-web-sm==3.7.1

Downloading https://github.com/explosion/spacy-models/releases/download/en_core_web_12.8/12.8 MB 103.4 MB/s eta 0:00:00

```
Requirement already satisfied: spacy<3.8.0,>=3.7.2 in /usr/local/lib/python3.11/dist-
Requirement already satisfied: spacy-legacy<3.1.0,>=3.0.11 in /usr/local/lib/python3.
Requirement already satisfied: spacy-loggers<2.0.0,>=1.0.0 in /usr/local/lib/python3.
Requirement already satisfied: murmurhash<1.1.0,>=0.28.0 in /usr/local/lib/python3.11
Requirement already satisfied: cymem<2.1.0,>=2.0.2 in /usr/local/lib/python3.11/dist-
Requirement already satisfied: preshed<3.1.0,>=3.0.2 in /usr/local/lib/python3.11/dis-
Requirement already satisfied: thinc<8.3.0,>=8.2.2 in /usr/local/lib/python3.11/dist-
Requirement already satisfied: wasabi<1.2.0,>=0.9.1 in /usr/local/lib/python3.11/dist-
Requirement already satisfied: srsly<3.0.0,>=2.4.3 in /usr/local/lib/python3.11/dist-
Requirement already satisfied: catalogue<2.1.0,>=2.0.6 in /usr/local/lib/python3.11/c
Requirement already satisfied: weasel<0.5.0,>=0.1.0 in /usr/local/lib/python3.11/dist-
Requirement already satisfied: typer<1.0.0,>=0.3.0 in /usr/local/lib/python3.11/dist-
Requirement already satisfied: tomviz <-1> in /usr/local/lib/python3.11/dist-
```

In this step, I set up the NLP preprocessing pipeline by installing and importing key libraries like SpaCy and NLTK. I downloaded the required stopwords and loaded SpaCy's language model, which I used later for tokenization and text cleaning.

```
1 # Define text preprocessing function
2 def clean_text(text):
3     # Remove special characters, digits, and punctuation
4     text = re.sub(r'\d+', '', text) # Remove digits
5     text = text.translate(str.maketrans('', '', string.punctuation)) # Remove punctuation
6     text = text.lower() # Convert text to lowercase
7
8     # Remove stopwords
9     stop_words = set(stopwords.words('english'))
10    words = text.split()
11    filtered_words = [word for word in words if word not in stop_words]
12
13    return " ".join(filtered_words)
14
15

1 # Apply text cleaning to 'Full Article Content' column
2 df["Cleaned Text"] = df["Full Article Content"].apply(clean_text)
3
4 # Tokenization using SpaCy
5 df["Tokens"] = df["Cleaned Text"].apply(lambda text: [token.text for token in nlp(text)])
6
7 # Lemmatization using SpaCy
8 df["Lemmatized Text"] = df["Cleaned Text"].apply(lambda text: " ".join([token.lemma_ for
9
10
```

In this step, I cleaned the text data by removing special characters, digits, punctuation, and stopwords. Then, I used Spacy to tokenize the text and applied lemmatization to reduce words to their root form. This process made the text clean and structured, preparing it for further analysis like sentiment detection and topic modeling.

▼ Step 6: Convert Text Data Using TF-IDF

```

1 from sklearn.feature_extraction.text import TfidfVectorizer
2
3 # Initialize TF-IDF Vectorizer
4 tfidf_vectorizer = TfidfVectorizer(max_features=500, stop_words="english") # Top 500 fe
5
6 # Apply TF-IDF transformation
7 tfidf_matrix = tfidf_vectorizer.fit_transform(df["Cleaned Text"])
8
9 # Convert TF-IDF matrix to a DataFrame
10 tfidf_df = pd.DataFrame(tfidf_matrix.toarray(), columns=tfidf_vectorizer.get_feature_nam
11
12 # Combine with original dataset
13 df_tfidf = pd.concat([df, tfidf_df], axis=1)
14
15 # Save the final dataset
16 df_tfidf.to_csv("nytimes_articles_tfidf.csv", index=False)
17
18 # Display the transformed dataset
19 print("TF-IDF transformation completed! Here is a preview:")
20 print(df_tfidf.head())
21

```

→ TF-IDF transformation completed! Here is a preview:

Title \

0	She Fell in Love With ChatGPT. Like, Actual Lo...
1	ChatGPT Is Restored After Hourslong Outage
2	She Is in Love With ChatGPT
3	How Helpful Is Operator, OpenAI's New A.I. Agent?
4	They Invested Billions. Then the A.I. Script G...

Publication Date \

0	2025-02-25 11:00:10+00:00
1	2025-01-23 13:02:14+00:00
2	2025-01-15 14:00:09+00:00
3	2025-02-01 10:01:31+00:00
4	2025-01-29 10:03:18+00:00

Author Location \

0	By Natalie Kitroeff, Kashmir Hill, Nina Feldma...	Podcasts
1	By Victor Mather	Business Day
2	By Kashmir Hill	Technology

```
3
4
```

By Kevin Roose	Technology
By Erin Griffith	Technology

Summary \

- 0 The story of a woman who is dating an A.I. cha...
- 1 Thousands of users were unable to access OpenA...
- 2 A 28-year-old woman with a busy social life sp...
- 3 Operator, a new computer-using tool from OpenA...
- 4 Venture capitalists plowed money into A.I. sta...

Full Article Content \

- 0 Warning: This episode discusses sexual themes.
- 1 ChatGPT was disrupted for several hours on Thu...
- 2 Ayrin's love affair with her A.I. boyfriend st...
- 3 In the past week, OpenAI's Operator has done t...
- 4 Jordan Jacobs, an investor at the venture capi...

URL Word Count \

	URL	Word Count	\
0	https://www.nytimes.com/2025/02/25/podcasts/the-ayrin-operator.html	6	
1	https://www.nytimes.com/2025/01/23/business/chatgpt-disruption.html	22	
2	https://www.nytimes.com/2025/01/15/technology/ayrin-operator.html	10	
3	https://www.nytimes.com/2025/02/01/technology/ayrin-operator.html	13	
4	https://www.nytimes.com/2025/01/29/technology/ayrin-operator.html	47	

	Sentence Count	Character Count	...	world	worried	worse	wowed	write	\
0	2	46	...	0.0	0.0	0.0	0.0	0.0	
1	2	156	...	0.0	0.0	0.0	0.0	0.0	
2	4	64	...	0.0	0.0	0.0	0.0	0.0	
3	1	73	...	0.0	0.0	0.0	0.0	0.0	
4	3	282	...	0.0	0.0	0.0	0.0	0.0	

	writer	writing	year	years	york
0	0.0	0.0	0.0	0.0	0.0
1	0.0	0.0	0.0	0.0	0.0
2	0.0	0.0	0.0	0.0	0.0
3	0.0	0.0	0.0	0.0	0.0
4	0.0	0.0	0.0	0.0	0.0

-- -- -- -- --

In this step, I converted the cleaned text into numerical features using the TF-IDF vectorizer. I limited it to the top 500 important terms to focus on the most relevant words. This transformation made the text data ready for machine learning models and further analysis.

▼ Step 7: NLP Application 1 – Sentiment Analysis

```
1 # Install VADER for sentiment analysis
2 !pip install vaderSentiment
3
4 # Import necessary libraries
5 import pandas as pd
```

```
6 import matplotlib.pyplot as plt
7 from vaderSentiment.vaderSentiment import SentimentIntensityAnalyzer
8
9
10 # Initialize VADER sentiment analyzer
11 analyzer = SentimentIntensityAnalyzer()
12
13
```

→ Collecting vaderSentiment

```
  Downloading vaderSentiment-3.3.2-py2.py3-none-any.whl.metadata (572 bytes)
Requirement already satisfied: requests in /usr/local/lib/python3.11/dist-packages (from vaderSentiment)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.11/dist-packages
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.11/dist-packages
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.11/dist-packages
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.11/dist-packages
  Downloading vaderSentiment-3.3.2-py2.py3-none-any.whl (125 kB)
                                         126.0/126.0 kB 11.0 MB/s eta 0:00:00
```

```
Installing collected packages: vaderSentiment
```

```
Successfully installed vaderSentiment-3.3.2
```



```
1 # Function to get sentiment scores
2 def get_sentiment(text):
3     sentiment = analyzer.polarity_scores(str(text)) # Ensure text is string
4     return sentiment["compound"] # Use compound score for overall sentiment
5
6 # Apply sentiment analysis
7 df["Sentiment Score"] = df["Cleaned Text"].apply(get_sentiment)
8
9 # Categorize sentiment (Positive, Neutral, Negative)
10 df["Sentiment Category"] = df["Sentiment Score"].apply(lambda x: "Positive" if x > 0.05
11
12 # Save sentiment analysis results
13 df.to_csv("nytimes_articles_sentiment.csv", index=False)
14
15 # Plot sentiment trends over time
16 df["Publication Date"] = pd.to_datetime(df["Publication Date"], errors="coerce")
17 df["Year-Month"] = df["Publication Date"].dt.to_period("M").astype(str)
18
19 # Aggregate sentiment scores per month
20 sentiment_trend = df.groupby("Year-Month")["Sentiment Score"].mean()
21
22
```

→ <ipython-input-20-5026b2f8f2c3>:17: UserWarning:

```
Converting to PeriodArray/Index representation will drop timezone information.
```

I used VADER sentiment analysis to evaluate the tone of each article, assigning a sentiment score and category (Positive, Negative, or Neutral). Then, I analyzed how sentiment trends changed over time by aggregating the results monthly. This helped identify shifts in media sentiment on the selected topics.

```
1 sentiment_trend_df = sentiment_trend.reset_index()
2 sentiment_trend_df.columns = ["Year-Month", "Average Sentiment Score"]
3
4 # Convert 'Year-Month' Period to string for Plotly compatibility
5 sentiment_trend_df["Year-Month"] = sentiment_trend_df["Year-Month"].astype(str)
6
7 # Create line chart using Plotly
8 fig = px.line(
9     sentiment_trend_df,
10    x="Year-Month",
11    y="Average Sentiment Score",
12    title="Sentiment Trends Over Time",
13    markers=True
14 )
15
16 # Add gridlines and customize layout
17 fig.update_layout(
18     xaxis_title="Year-Month",
19     yaxis_title="Average Sentiment Score",
20     plot_bgcolor='white',
21     width=900,
22     height=500,
23     xaxis=dict(
24         showgrid=True,           # Enable X-axis gridlines
25         gridcolor='lightgrey',
26         tickangle=-45
27     ),
28     yaxis=dict(
29         showgrid=True,           # Enable Y-axis gridlines
30         gridcolor='lightgrey'
31     )
32 )
33
34 # Show the plot
35 fig.show()
36
```



Visualized the sentiment trends over time using a Plotly line chart. The graph shows how the average sentiment score of articles changed month by month, with added gridlines for clarity and easier interpretation.

```
1 # Display final dataset preview
2 print("Sentiment analysis completed! Here is a preview:")
3 print(df[["Title", "Sentiment Score", "Sentiment Category"]].head())
4
```

Sentiment analysis completed! Here is a preview:

	Title	Sentiment Score	\
0	She Fell in Love With ChatGPT. Like, Actual Lo...	-0.3400	
1	ChatGPT Is Restored After Hourslong Outage	0.0000	
2	She Is in Love With ChatGPT	0.6369	
3	How Helpful Is Operator, OpenAI's New A.I. Agent?	0.0000	
4	They Invested Billions. Then the A.I. Script G...	0.4767	

	Sentiment Category
0	Negative
1	Neutral
2	Positive
3	Neutral
4	Positive

✓ Step 8: Word Cloud Analysis

```
1 # Install wordcloud if not already installed
2 !pip install wordcloud
3
4 # Import necessary libraries
5 import pandas as pd
6 import matplotlib.pyplot as plt
7 from wordcloud import WordCloud
8
9 # Load the sentiment analysis dataset
10 csv_filepath = "nytimes_articles_sentiment.csv" # Ensure correct file path
11 df = pd.read_csv(csv_filepath)
12
13 # Separate positive and negative articles
14 positive_text = " ".join(df[df["Sentiment Category"] == "Positive"]["Cleaned Text"].dropna())
15 negative_text = " ".join(df[df["Sentiment Category"] == "Negative"]["Cleaned Text"].dropna())
16
```

→ Requirement already satisfied: wordcloud in /usr/local/lib/python3.11/dist-packages (1.9)
Requirement already satisfied: numpy>=1.6.1 in /usr/local/lib/python3.11/dist-packages (from wordcloud)
Requirement already satisfied: pillow in /usr/local/lib/python3.11/dist-packages (from wordcloud)
Requirement already satisfied: matplotlib in /usr/local/lib/python3.11/dist-packages (from wordcloud)
Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.11/dist-packages (from wordcloud)
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.11/dist-packages (from wordcloud)
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.11/dist-packages (from wordcloud)
Requirement already satisfied: kiwisolver>=1.3.1 in /usr/local/lib/python3.11/dist-packages (from wordcloud)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.11/dist-packages (from wordcloud)
Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.11/dist-packages (from wordcloud)
Requirement already satisfied: python-dateutil>=2.7 in /usr/local/lib/python3.11/dist-packages (from wordcloud)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.11/dist-packages (from wordcloud)

After completing sentiment analysis, I reviewed the results by displaying a preview of article titles along with their sentiment scores and categories. Then, I prepared for word cloud analysis by separating positive and negative texts, which will help visualize the most common words used in articles with different sentiments.

Positive Reviews

```
1 # Generate Word Cloud for Positive Words
2 plt.figure(figsize=(10, 5))
3 wordcloud_positive = WordCloud(width=800, height=400, background_color="white").generate(positive_text)
4 plt.imshow(wordcloud_positive, interpolation="bilinear")
5 plt.axis("off")
```

```
6 plt.title("Word Cloud for Positive Sentiment Words")
7 plt.show()
8
9
```



Negative` Reviews

```
1 # Generate Word Cloud for Negative Words
2 plt.figure(figsize=(10, 5))
3 wordcloud_negative = WordCloud(width=800, height=400, background_color="black", colormap="red")
4 plt.imshow(wordcloud_negative, interpolation="bilinear")
5 plt.axis("off")
6 plt.title("Word Cloud for Negative Sentiment Words")
7 plt.show()
8
```



- I generated two separate word clouds to visualize the most frequently used words in positive and negative sentiment articles.
- For positive sentiment, the word cloud highlighted terms like ChatGPT, artificial intelligence, and technology, often associated with innovation and progress.
- For negative sentiment, common terms included chatbot, openai, and google, often appearing in contexts expressing concern or criticism.
- The word clouds helped quickly identify dominant themes and language patterns in articles with differing sentiments.

▼ Step 9: NLP Application 2 – Topic Modeling

```
1 # Install necessary libraries if not already installed
2 !pip install gensim
3
4 # Import necessary libraries
5 import pandas as pd
6 import gensim
7 from gensim import corpora
8 from gensim.models import LdaModel
9 from nltk.tokenize import word_tokenize
10 import nltk
```

```

11 import matplotlib.pyplot as plt
12
13 # Ensure NLTK tokenization resources are available
14 nltk.download("punkt")
15
16 # Load the cleaned dataset
17 csv_filepath = "/content/nytimes_articles_sentiment.csv"
18
19

```

→ Requirement already satisfied: gensim in /usr/local/lib/python3.11/dist-packages (4.3.3)
 Requirement already satisfied: numpy<2.0,>=1.18.5 in /usr/local/lib/python3.11/dist-pac
 Requirement already satisfied: scipy<1.14.0,>=1.7.0 in /usr/local/lib/python3.11/dist-pa
 Requirement already satisfied: smart-open>=1.8.1 in /usr/local/lib/python3.11/dist-pac
 Requirement already satisfied: wrapt in /usr/local/lib/python3.11/dist-packages (from sn
 [nltk_data] Downloading package punkt to /root/nltk_data...
 [nltk_data] Unzipping tokenizers/punkt.zip.

```

1 import nltk
2 nltk.download('punkt')
3
4 import nltk
5 nltk.download('punkt_tab')
6
7 # Tokenize the cleaned text data
8 df["Tokenized Text"] = df["Cleaned Text"].apply(lambda x: word_tokenize(str(x)))
9
10

```

→ [nltk_data] Downloading package punkt to /root/nltk_data...
 [nltk_data] Package punkt is already up-to-date!
 [nltk_data] Downloading package punkt_tab to /root/nltk_data...
 [nltk_data] Unzipping tokenizers/punkt_tab.zip.

```

1 # Create dictionary and corpus for LDA
2 dictionary = corpora.Dictionary(df["Tokenized Text"])
3 corpus = [dictionary.doc2bow(text) for text in df["Tokenized Text"]]
4
5 # Train LDA model (adjust num_topics as needed)
6 num_topics = 5 # Define number of topics to extract
7 lda_model = LdaModel(corpus, num_topics=num_topics, id2word=dictionary, passes=10)
8
9

```

I used the Gensim library to apply Latent Dirichlet Allocation (LDA) for topic modeling on the cleaned text data. After tokenizing the articles using NLTK, I created a dictionary and corpus, which served as inputs for the LDA model. The model was trained to identify the top 5 key topics across the articles, helping uncover the main themes present in the dataset.

```

1 # Display the top words for each topic
2 topics = lda_model.print_topics(num_words=10)
3 print("\nExtracted Topics:")
4 for topic in topics:
5     print(topic)
6
7 # Assign topics to articles
8 df["Dominant Topic"] = [max(lda_model[doc], key=lambda x: x[1])[0] if lda_model[doc] else
9
10 # Analyze topic trends over time
11 df["Publication Date"] = pd.to_datetime(df["Publication Date"], errors="coerce")
12 df["Year-Month"] = df["Publication Date"].dt.to_period("M")
13
14 # Aggregate topic trends
15 topic_trends = df.groupby(["Year-Month", "Dominant Topic"]).size().unstack(fill_value=0)
16
17

```



Extracted Topics:

```

(0, '0.012*"chatbot" + 0.012*"chatgpt" + 0.010*"artificial" + 0.010*"openai" + 0.008*"ir
(1, '0.018*'" + 0.015*"s" + 0.012*"openai" + 0.011*"chatgpt" + 0.008*"chatbot" + 0.008*
(2, '0.036*'" + 0.022*"s" + 0.010*"- + 0.010*"like" + 0.010*"ai" + 0.008*"" + 0.008*"
(3, '0.018*"openai" + 0.016*"chatgpt" + 0.014*'" + 0.013*"intelligence" + 0.013*"artifi
(4, '0.020*'" + 0.012*"chatgpt" + 0.011*"read" + 0.011*"artificial" + 0.011*"i" + 0.011
<ipython-input-29-7ccb3929e74d>:12: UserWarning:

```

Converting to PeriodArray/Index representation will drop timezone information.



I displayed the top words for each topic generated by the LDA model to understand the main themes. Then, I assigned a dominant topic to each article and analyzed how these topics trended over time by grouping them monthly. This helped visualize shifts in focus across different periods.

```

1 # Reset index if necessary and convert Year-Month to string for Plotly
2 topic_trends_df = topic_trends.reset_index()
3 topic_trends_df["Year-Month"] = topic_trends_df["Year-Month"].astype(str)
4
5 # Melt the DataFrame to long format for Plotly
6 topic_trends_melted = topic_trends_df.melt(id_vars="Year-Month",
7                                              var_name="Topic",
8                                              value_name="Number of Articles")
9
10 # Create interactive line plot
11 fig = px.line(
12     topic_trends_melted,
13     x="Year-Month",

```

```
14     y="Number of Articles",
15     color="Topic",
16     markers=True,
17     title="Topic Trends Over Time"
18 )
19
20 # Customize layout with grids
21 fig.update_layout(
22     xaxis_title="Year-Month",
23     yaxis_title="Number of Articles",
24     legend_title="Topics",
25     width=950,
26     height=600,
27     plot_bgcolor='white',
28     xaxis=dict(
29         showgrid=True,           # Enable grid lines for X-axis
30         gridcolor='lightgrey',
31         tickangle=-45
32     ),
33     yaxis=dict(
34         showgrid=True,           # Enable grid lines for Y-axis
35         gridcolor='lightgrey'
36     )
37 )
38
39 # Show the plot
40 fig.show()
41
```



I visualized how different topics trended over time by creating an interactive line chart using Plotly. The graph shows the number of articles for each topic per month, with customized grid lines for better readability. This helped track shifts in focus across topics throughout the analyzed period.

```

1 # Save topic modeling results
2 df.to_csv("nytimes_articles_topics.csv", index=False)
3
4 # Display preview of dataset with topics
5 print("\n Topic modeling completed! Preview of dataset with assigned topics:")
6 print(df[["Title", "Dominant Topic"]].head())
7

```



Topic modeling completed! Preview of dataset with assigned topics:

	Title	Dominant Topic
0	She Fell in Love With ChatGPT. Like, Actual Lo...	4
1	ChatGPT Is Restored After Hourslong Outage	0
2	She Is in Love With ChatGPT	2
3	How Helpful Is Operator, OpenAI's New A.I. Agent?	1
4	They Invested Billions. Then the A.I. Script G...	4

Workflow - Building a custom Sentiment Classifier

