

Team Members

Name	SRN
Naren Chandrashekhar	PES2UG20CS216
Nitin Jayachandran	PES2UG20CS231
Rohan Amin	PES2UG20CS416
Anurag B	PES2UG20CS404

Section 1

- Corpus Details: The BBC News Dataset contains 2225 news articles from the BBC, spanning 5 different categories: business, entertainment, politics, sport, and tech. Each article is labeled with its corresponding category.

The dataset is stored in a CSV file, with each row representing a single article and the following columns:

category: the category label for the article (one of "business", "entertainment", "politics", "sport", or "tech").

text: the full text of the article.

- Source: <https://www.kaggle.com/datasets/sahilkirpekar/bbcnews-dataset?resource=download> (<https://www.kaggle.com/datasets/sahilkirpekar/bbcnews-dataset?resource=download>)

Libraries Used

```
In [10]: import pandas as pd
import re
import nltk
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
from nltk.tokenize import word_tokenize
nltk.download('stopwords')
nltk.download('punkt')
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics.pairwise import cosine_similarity

from collections import defaultdict
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
```

Knowing the Dataset

```
In [11]: df = pd.read_csv("/content/BBCNews.csv")
print(df.head())
```

```
Unnamed: 0                               descr \
0          0  chelsea sack mutu  chelsea have sacked adrian ...
1          1  record fails to lift lacklustre meet  yelena i...
2          2  edu describes tunnel fracas  arsenals edu has ...
3          3  ogara revels in ireland victory  ireland flyha...
4          4  unclear future for striker baros  liverpool fo...

tags
0  sports, stamford bridge, football association, ...
1  sports, madrid, birmingham, france, scotland, ...
2  sports, derby, brazil, tunnel fracasedu, food, ...
3  sports, bbc, united kingdom, ireland, brian o'...
4  sports, liverpool, daily sport, millennium sta...
```

Preprocessing Of Data

```
In [12]: #Lowercase all the words in the dataset
df['descr'] = df['descr'].str.lower()

# display the updated dataset
print(df.head())
```

```
Unnamed: 0                               descr \
0      0 chelsea sack mutu chelsea have sacked adrian ...
1      1 record fails to lift lacklustre meet yelena i...
2      2 edu describes tunnel fracas arsenals edu has ...
3      3 ogara revels in ireland victory ireland flyha...
4      4 unclear future for striker baros liverpool fo...

                                         tags
0  sports, stamford bridge, football association, ...
1  sports, madrid, birmingham, france, scotland, ...
2  sports, derby, brazil, tunnel fracasedu, food, ...
3  sports, bbc, united kingdom, ireland, brian o'...
4  sports, liverpool, daily sport, millennium sta...
```

```
In [13]: #Remove stop words
# create a set of stop words
stop_words = set(stopwords.words('english'))

# remove stop words from the 'descr' column
df['descr'] = df['descr'].apply(lambda x: ' '.join([word for word in x.split() if word not in stop_words]))

# display the updated dataset
print(df.head())
```

```
Unnamed: 0                               descr \
0      0 chelsea sack mutu chelsea sacked adrian mutu f...
1      1 record fails lift lacklustre meet yelena isinb...
2      2 edu describes tunnel fracas arsenals edu lifte...
3      3 ogara revels ireland victory ireland flyhalf r...
4      4 unclear future striker baros liverpool forward...

                                         tags
0  sports, stamford bridge, football association, ...
1  sports, madrid, birmingham, france, scotland, ...
2  sports, derby, brazil, tunnel fracasedu, food, ...
3  sports, bbc, united kingdom, ireland, brian o'...
4  sports, liverpool, daily sport, millennium sta...
```

```
In [14]: # Initialize the stemmer
stemmer = PorterStemmer()

# Define a function to stem the text
def stem_text(text):
    # Tokenize the text
    tokens = word_tokenize(text)
    # Stem each token and add to a new list
    stemmed_tokens = [stemmer.stem(token) for token in tokens]
    # Join the tokens back into a string
    text = ' '.join(stemmed_tokens)
    return text

# Apply the stem function to the 'clean_text' column of the dataframe
df['descr'] = df['descr'].apply(stem_text)
print(df.head())
```

```
Unnamed: 0                               descr \
0          0  chelsea sack mutu chelsea sack adrian mutu fai...
1          1  record fail lift lacklustr meet yelena isinbay...
2          2  edu describ tunnel fraca arsen edu lift lid sc...
3          3  ogara revel ireland victori ireland flyhalf ro...
4          4  unclear futur striker baro liverpool forward m...

tags
0  sports, stamford bridge, football association, ...
1  sports, madrid, birmingham, france, scotland, ...
2  sports, derby, brazil, tunnel fracasedu, food, ...
3  sports, bbc, united kingdom, ireland, brian o'...
4  sports, liverpool, daily sport, millennium sta...
```

```
In [15]: # Define a function to remove special characters
def remove_special_chars(text):
    # Define the pattern to match special characters
    pattern = r'[^a-zA-Z0-9\s]'
    # Replace special characters with a space
    text = re.sub(pattern, ' ', text)
    return text

# Apply the function to the 'clean_text' column of the dataframe
df['descr'] = df['descr'].apply(remove_special_chars)
print(df.head())
```

	Unnamed: 0	descr \
0	0	chelsea sack mutu chelsea sack adrian mutu fai...
1	1	record fail lift lacklustr meet yelena isinbay...
2	2	edu describ tunnel fraca arsen edu lift lid sc...
3	3	ogara revel ireland victori ireland flyhalf ro...
4	4	unclear futur striker baro liverpool forward m...

	tags
0	sports, stamford bridge, football association,...
1	sports, madrid, birmingham, france, scotland, ...
2	sports, derby, brazil, tunnel fracasedu, food,...
3	sports, bbc, united kingdom, ireland, brian o'...
4	sports, liverpool, daily sport, millennium sta...

Section 2

Set of

- Free Text Test Queries
- Wild card Queries
- Phase Queries

Free Text Test Queries

- For Boolean Retrieval: Results shown in Section 4.1

1. query: white house terrorist
2. query: federer championship
3. query: ronaldo manchester united
4. query: tax law fraud
5. query: windows microsoft
6. query: economy us bank loan

- For Inverted Index: Results shown in Section 4.2

1. query: UK election
2. query: white house
3. query: chelsea
4. query: robot

5. query: hollywood
6. query: apple
7. query: blockchain
8. query: pear

Wild Card Queries

Results shown in Section 4.3

1. query: robo*
2. query: *tech
3. query: cyber*attack
4. query: mark*crash
5. query: energy

Phrase Queries

Results shown in Section 4.4

1. query: climate change
2. query: gender pay gap
3. query: global economic
4. query: food waste reduction initiatives
5. query: human rights
6. query: economic growth

Section 3

Data structures used with brief reasons and similarity scheme

Data Structures Used

1. Inverted Index: An inverted index is a data structure that stores a mapping between terms and the documents that contain them. Inverted index is efficient for keyword-based search because it can quickly locate documents containing a particular term or a combination of terms. It is widely used in search engines to efficiently retrieve relevant documents based on user queries.
2. Sets: A set is an unordered collection of unique items, which can be used to store and manipulate sets of terms that occur in a document or query. In particular, sets can be used to represent the set of terms that appear in a document or query, or to store the set of documents that contain a particular term. One common use of sets in information retrieval is for document representation using the bag-of-words model. In this model, a document is represented as a set of terms, where the frequency of each term in the document is not considered. The set of terms that appear in a document can be efficiently computed using a set data structure. Another use of sets in information retrieval is for computing set operations such as union, intersection, and difference.
3. Dictionaries: A dictionary is a collection of key-value pairs, where each key is unique and maps to a corresponding value. In information retrieval, dictionaries can be used to represent mappings between terms and their corresponding document frequency or term frequency in a collection of documents. One common use of dictionaries in information

retrieval is for building an inverted index. In this index, each term in the vocabulary is mapped to a list of document IDs that contain that term. The inverted index can be implemented using a dictionary, where each key is a term and each value is a list of document IDs. The dictionary can be efficiently updated as new documents are added to the collection, and can be used to quickly retrieve the set of documents that contain a given term.

4. Lists: A list is an ordered collection of items, which can be used to store and manipulate sequences of terms or documents in a particular order. In particular, lists can be used to represent sequences of documents that match a query, or to store sequences of terms in a document. One common use of lists in information retrieval is for document ranking using the vector space model. In this model, each document is represented as a vector of weights, where each weight corresponds to the importance of a particular term in the document. To rank documents based on their similarity to a query, the model computes the cosine similarity between the query vector and each document vector, and returns a ranked list of documents. The ranked list of documents can be represented as a list of document IDs in order of decreasing similarity score.

Similarity Scheme

Similarity schemes are methods used to measure the similarity between a query and a document in information retrieval. The goal of a similarity scheme is to quantify how relevant a document is to a query based on the similarity between their respective representations.

1. Cosine similarity: This is a commonly used similarity scheme in vector space models. Cosine similarity measures the cosine of the angle between the query vector and the document vector, and is given by the dot product of the two vectors divided by their magnitudes. Cosine similarity is widely used in text retrieval because it is efficient and easy to implement.
2. BM25: This is a similarity scheme that is commonly used in ranking algorithms for text retrieval. BM25 measures the relevance of a document to a query based on the frequency of the query terms in the document and the frequency of the query terms in the entire collection of documents. It is designed to be robust to differences in document lengths and to handle rare terms.
3. TF-IDF (Term Frequency-Inverse Document Frequency) is a weighting scheme commonly used in information retrieval to rank documents based on their relevance to a query. While TF-IDF is not a similarity scheme in and of itself, it is often used as a component of similarity schemes such as the vector space model. The basic idea behind TF-IDF is to weigh the importance of a term in a document based on how frequently it appears in the document (term frequency) and how rarely it appears in the collection of all documents (inverse document frequency). This is done by multiplying the term frequency by the inverse document frequency, resulting in a weight that is high if the term appears frequently in the document and rarely in the collection, and low if the term appears rarely in the document or frequently in the collection.

Section 4.1

Result of Boolean Retrieval: On free text queries

Below are some Boolean Retrieved documents which specify the document id where the queries have found a match. The documents contain each of the given words in the queries. We have provided some useful queries for which the required documents have been specified

```
In [16]: df = pd.read_csv("/content/BBCNews.csv")
# define the boolean retrieval function
def boolean_search(query, df):
    # split the query into individual terms
    query_terms = query.lower().split()

    # initialize an empty set of documents
    result = set()

    # Loop over the documents in the dataframe
    for i in range(len(df)):
        # get the text of the current document
        text = df.iloc[i]["descr"].lower()

        # check if the document contains all of the query terms
        contains_all_terms = all(term in text for term in query_terms)

        # if the document contains all of the query terms, add it to the result
        if contains_all_terms:
            result.add(i)

    return result

# test the boolean retrieval function
query = "white house terrorist"
results = boolean_search(query, df)
print(results)
```

```
{1186}
```

```
In [17]: df = pd.read_csv("/content/BBCNews.csv")
# define the boolean retrieval function
def boolean_search(query, df):
    # split the query into individual terms
    query_terms = query.lower().split()

    # initialize an empty set of documents
    result = set()

    # Loop over the documents in the dataframe
    for i in range(len(df)):
        # get the text of the current document
        text = df.iloc[i]["descr"].lower()

        # check if the document contains all of the query terms
        contains_all_terms = all(term in text for term in query_terms)

        # if the document contains all of the query terms, add it to the result
        if contains_all_terms:
            result.add(i)

    return result

# test the boolean retrieval function
query = "federer championship"
results = boolean_search(query, df)
print(results)
```

```
{339, 157, 31}
```

```
In [18]: df = pd.read_csv("/content/BBCNews.csv")
# define the boolean retrieval function
def boolean_search(query, df):
    # split the query into individual terms
    query_terms = query.lower().split()

    # initialize an empty set of documents
    result = set()

    # Loop over the documents in the dataframe
    for i in range(len(df)):
        # get the text of the current document
        text = df.iloc[i]["descr"].lower()

        # check if the document contains all of the query terms
        contains_all_terms = all(term in text for term in query_terms)

        # if the document contains all of the query terms, add it to the result
        if contains_all_terms:
            result.add(i)

    return result

# test the boolean retrieval function
query = "ronaldo manchester united"
results = boolean_search(query, df)
print(results)
```

```
{192, 5, 7, 1551, 1491, 116}
```

```
In [19]: df = pd.read_csv("/content/BBCNews.csv")
# define the boolean retrieval function
def boolean_search(query, df):
    # split the query into individual terms
    query_terms = query.lower().split()

    # initialize an empty set of documents
    result = set()

    # Loop over the documents in the dataframe
    for i in range(len(df)):
        # get the text of the current document
        text = df.iloc[i]["descr"].lower()

        # check if the document contains all of the query terms
        contains_all_terms = all(term in text for term in query_terms)

        # if the document contains all of the query terms, add it to the result
        if contains_all_terms:
            result.add(i)

    return result

# test the boolean retrieval function
query = "tax law fraud"
results = boolean_search(query, df)
print(results)
```

{1729, 1763, 1891, 1915, 2276, 2373, 1737, 1900, 2235, 1679, 2389, 1656, 1177, 1659, 1630}

```
In [20]: df = pd.read_csv("/content/BBCNews.csv")
# define the boolean retrieval function
def boolean_search(query, df):
    # split the query into individual terms
    query_terms = query.lower().split()

    # initialize an empty set of documents
    result = set()

    # Loop over the documents in the dataframe
    for i in range(len(df)):
        # get the text of the current document
        text = df.iloc[i]["descr"].lower()

        # check if the document contains all of the query terms
        contains_all_terms = all(term in text for term in query_terms)

        # if the document contains all of the query terms, add it to the result
        if contains_all_terms:
            result.add(i)

    return result

# test the boolean retrieval function
query = "windows microsoft"
results = boolean_search(query, df)
print(results)
```

```
{1272, 1282, 1413, 774, 778, 1290, 1418, 1293, 1294, 408, 413, 544, 800, 803, 676, 677, 678, 806, 810, 434, 690, 567, 440, 573, 701, 449, 581, 454, 712, 459, 464, 1362, 595, 1237, 472, 475, 476, 604, 1252, 1385, 1258, 620, 621, 1391, 1267, 756, 1399, 760, 1277}
```

```
In [21]: df = pd.read_csv("/content/BBCNews.csv")
# define the boolean retrieval function
def boolean_search(query, df):
    # split the query into individual terms
    query_terms = query.lower().split()

    # initialize an empty set of documents
    result = set()

    # Loop over the documents in the dataframe
    for i in range(len(df)):
        # get the text of the current document
        text = df.iloc[i]["descr"].lower()

        # check if the document contains all of the query terms
        contains_all_terms = all(term in text for term in query_terms)

        # if the document contains all of the query terms, add it to the result
        if contains_all_terms:
            result.add(i)

    return result

# test the boolean retrieval function
query = "economy us loan bank"
results = boolean_search(query, df)
print(results)
```

{2249, 1771, 1709, 2349, 1746, 1714, 2228}

Section 4.2

Result with inverted index: On free text queries with rank.

The first cell block is a basic code to give the inverted index for each word in the document. The following cell blocks give the result of a free text query with the ranking.

Explanation as to how the query works, various examples with explanation has been provided after each cell output.

Ranking here specifies the most relevant document retrieved with respect to the query given. The output of the document ranking is in descending order.

```
In [22]: # Define a regular expression to tokenize the text
tokenizer = re.compile(r"\w+")

# Create an empty dictionary to store the inverted index
inverted_index = defaultdict(list)

# Loop over each document in the dataset
for i, row in df.iterrows():
    # Tokenize the text
    tokens = tokenizer.findall(row['descr'])
    # Count the frequency of each token
    term_freq = defaultdict(int)
    for token in tokens:
        term_freq[token] += 1
    # Add the document to the inverted index for each unique token
    for token, freq in term_freq.items():
        inverted_index[token].append((i, freq))

# Print the inverted index
for token, postings in inverted_index.items():
    print(f"{token}: {postings}")
```

```
chelsea: [(0, 8), (4, 1), (12, 8), (24, 3), (25, 1), (42, 2), (47, 2), (10
8, 1), (110, 5), (122, 1), (124, 1), (126, 4), (129, 1), (137, 2), (139,
8), (141, 3), (144, 2), (146, 2), (151, 1), (154, 3), (155, 8), (163, 4),
(166, 1), (192, 1), (197, 5), (346, 3), (358, 7), (361, 4), (362, 2), (36
5, 3), (383, 3), (389, 3), (392, 1), (395, 1), (396, 1), (453, 1), (631,
2), (632, 3), (635, 1), (637, 1), (638, 2), (643, 4), (652, 2), (662, 1),
(922, 1), (1187, 2), (1271, 1), (1451, 3), (1456, 1), (1462, 1), (1483,
3), (1485, 4), (1486, 5), (1493, 7), (1501, 2), (1509, 3), (1523, 2), (152
4, 3), (1534, 1), (1535, 1), (1540, 2), (1550, 1), (1558, 3), (1565, 7),
(1573, 3), (1581, 3), (1598, 4), (1604, 1), (2308, 3)]
sack: [(0, 2), (147, 1), (206, 1), (207, 1), (298, 1), (309, 1), (954, 1),
(1222, 1), (1427, 1), (1598, 2), (1605, 1)]
mutu: [(0, 6), (652, 6), (1598, 6)]
have: [(0, 5), (1, 10), (2, 1), (3, 5), (5, 1), (6, 1), (7, 1), (8, 1), (1
1, 2), (12, 3), (15, 1), (16, 3), (19, 1), (21, 3), (22, 4), (23, 2), (25,
2), (26, 7), (27, 6), (28, 1), (29, 3), (30, 6), (31, 1), (32, 3), (33,
3), (34, 6), (36, 1), (37, 1), (38, 1), (39, 1), (40, 5), (42, 1), (44,
1), (47, 8), (48, 1), (49, 3), (50, 2), (51, 1), (52, 1), (53, 3), (56,
4), (57, 2), (58, 1), (59, 3), (61, 2), (63, 2), (66, 2), (67, 3), (68,
1), (69, 1), (70, 1), (71, 1), (72, 1), (73, 1), (74, 1), (75, 1), (76, 1),
(77, 1), (78, 1), (79, 1), (80, 1), (81, 1), (82, 1), (83, 1), (84, 1), (85,
1), (86, 1), (87, 1), (88, 1), (89, 1), (90, 1), (91, 1), (92, 1), (93, 1),
(94, 1), (95, 1), (96, 1), (97, 1), (98, 1), (99, 1), (100, 1), (101, 1),
(102, 1), (103, 1), (104, 1), (105, 1), (106, 1), (107, 1), (108, 1), (109,
1), (110, 1), (111, 1), (112, 1), (113, 1), (114, 1), (115, 1), (116, 1),
(117, 1), (118, 1), (119, 1), (120, 1), (121, 1), (122, 1), (123, 1), (124,
1), (125, 1), (126, 1), (127, 1), (128, 1), (129, 1), (130, 1), (131, 1),
(132, 1), (133, 1), (134, 1), (135, 1), (136, 1), (137, 1), (138, 1), (139,
1), (140, 1), (141, 1), (142, 1), (143, 1), (144, 1), (145, 1), (146, 1),
(147, 1), (148, 1), (149, 1), (150, 1), (151, 1), (152, 1), (153, 1), (154,
1), (155, 1), (156, 1), (157, 1), (158, 1), (159, 1), (160, 1), (161, 1),
(162, 1), (163, 1), (164, 1), (165, 1), (166, 1), (167, 1), (168, 1), (169,
1), (170, 1), (171, 1), (172, 1), (173, 1), (174, 1), (175, 1), (176, 1),
(177, 1), (178, 1), (179, 1), (180, 1), (181, 1), (182, 1), (183, 1), (184,
1), (185, 1), (186, 1), (187, 1), (188, 1), (189, 1), (190, 1), (191, 1),
(192, 1), (193, 1), (194, 1), (195, 1), (196, 1), (197, 1), (198, 1), (199,
1), (200, 1), (201, 1), (202, 1), (203, 1), (204, 1), (205, 1), (206, 1),
(207, 1), (208, 1), (209, 1), (210, 1), (211, 1), (212, 1), (213, 1), (214,
1), (215, 1), (216, 1), (217, 1), (218, 1), (219, 1), (220, 1), (221, 1),
(222, 1), (223, 1), (224, 1), (225, 1), (226, 1), (227, 1), (228, 1), (229,
1), (230, 1), (231, 1), (232, 1), (233, 1), (234, 1), (235, 1), (236, 1),
(237, 1), (238, 1), (239, 1), (240, 1), (241, 1), (242, 1), (243, 1), (244,
1), (245, 1), (246, 1), (247, 1), (248, 1), (249, 1), (250, 1), (251, 1),
(252, 1), (253, 1), (254, 1), (255, 1), (256, 1), (257, 1), (258, 1), (259,
1), (260, 1), (261, 1), (262, 1), (263, 1), (264, 1), (265, 1), (266, 1),
(267, 1), (268, 1), (269, 1), (270, 1), (271, 1), (272, 1), (273, 1), (274,
1), (275, 1), (276, 1), (277, 1), (278, 1), (279, 1), (280, 1), (281, 1),
(282, 1), (283, 1), (284, 1), (285, 1), (286, 1), (287, 1), (288, 1), (289,
1), (290, 1), (291, 1), (292, 1), (293, 1), (294, 1), (295, 1), (296, 1),
(297, 1), (298, 1), (299, 1), (300, 1), (301, 1), (302, 1), (303, 1), (304,
1), (305, 1), (306, 1), (307, 1), (308, 1), (309, 1), (310, 1), (311, 1),
(312, 1), (313, 1), (314, 1), (315, 1), (316, 1), (317, 1), (318, 1), (319,
1), (320, 1), (321, 1), (322, 1), (323, 1), (324, 1), (325, 1), (326, 1),
(327, 1), (328, 1), (329, 1), (330, 1), (331, 1), (332, 1), (333, 1), (334,
1), (335, 1), (336, 1), (337, 1), (338, 1), (339, 1), (340, 1), (341, 1),
(342, 1), (343, 1), (344, 1), (345, 1), (346, 1), (347, 1), (348, 1), (349,
1), (350, 1), (351, 1), (352, 1), (353, 1), (354, 1), (355, 1), (356, 1),
(357, 1), (358, 1), (359, 1), (360, 1), (361, 1), (362, 1), (363, 1), (364,
1), (365, 1), (366, 1), (367, 1), (368, 1), (369, 1), (370, 1), (371, 1),
(372, 1), (373, 1), (374, 1), (375, 1), (376, 1), (377, 1), (378, 1), (379,
1), (380, 1), (381, 1), (382, 1), (383, 1), (384, 1), (385, 1), (386, 1),
(387, 1), (388, 1), (389, 1), (390, 1), (391, 1), (392, 1), (393, 1), (394,
1), (395, 1), (396, 1), (397, 1), (398, 1), (399, 1), (400, 1), (401, 1),
(402, 1), (403, 1), (404, 1), (405, 1), (406, 1), (407, 1), (408, 1), (409,
1), (410, 1), (411, 1), (412, 1), (413, 1), (414, 1), (415, 1), (416, 1),
(417, 1), (418, 1), (419, 1), (420, 1), (421, 1), (422, 1), (423, 1), (424,
1), (425, 1), (426, 1), (427, 1), (428, 1), (429, 1), (430, 1), (431, 1),
(432, 1), (433, 1), (434, 1), (435, 1), (436, 1), (437, 1), (438, 1), (439,
1), (440, 1), (441, 1), (442, 1), (443, 1), (444, 1), (445, 1), (446, 1),
(447, 1), (448, 1), (449, 1), (450, 1), (451, 1), (452, 1), (453, 1), (454,
1), (455, 1), (456, 1), (457, 1), (458, 1), (459, 1), (460, 1), (461, 1),
(462, 1), (463, 1), (464, 1), (465, 1), (466, 1), (467, 1), (468, 1), (469,
1), (470, 1), (471, 1), (472, 1), (473, 1), (474, 1), (475, 1), (476, 1),
(477, 1), (478, 1), (479, 1), (480, 1), (481, 1), (482, 1), (483, 1), (484,
1), (485, 1), (486, 1), (487, 1), (488, 1), (489, 1), (490, 1), (491, 1),
(492, 1), (493, 1), (494, 1), (495, 1), (496, 1), (497, 1), (498, 1), (499,
1), (500, 1), (501, 1), (502, 1), (503, 1), (504, 1), (505, 1), (506, 1),
(507, 1), (508, 1), (509, 1), (510, 1), (511, 1), (512, 1), (513, 1), (514,
1), (515, 1), (516, 1), (517, 1), (518, 1), (519, 1), (520, 1), (521, 1),
(522, 1), (523, 1), (524, 1), (525, 1), (526, 1), (527, 1), (528, 1), (529,
1), (530, 1), (531, 1), (532, 1), (533, 1), (534, 1), (535, 1), (536, 1),
(537, 1), (538, 1), (539, 1), (540, 1), (541, 1), (542, 1), (543, 1), (544,
1), (545, 1), (546, 1), (547, 1), (548, 1), (549, 1), (550, 1), (551, 1),
(552, 1), (553, 1), (554, 1), (555, 1), (556, 1), (557, 1), (558, 1), (559,
1), (560, 1), (561, 1), (562, 1), (563, 1), (564, 1), (565, 1), (566, 1),
(567, 1), (568, 1), (569, 1), (570, 1), (571, 1), (572, 1), (573, 1), (574,
1), (575, 1), (576, 1), (577, 1), (578, 1), (579, 1), (580, 1), (581, 1),
(582, 1), (583, 1), (584, 1), (585, 1), (586, 1), (587, 1), (588, 1), (589,
1), (590, 1), (591, 1), (592, 1), (593, 1), (594, 1), (595, 1), (596, 1),
(597, 1), (598, 1), (599, 1), (600, 1), (601, 1), (602, 1), (603, 1), (604,
1), (605, 1), (606, 1), (607, 1), (608, 1), (609, 1), (610, 1), (611, 1),
(612, 1), (613, 1), (614, 1), (615, 1), (616, 1), (617, 1), (618, 1), (619,
1), (620, 1), (621, 1), (622, 1), (623, 1), (624, 1), (625, 1), (626, 1),
(627, 1), (628, 1), (629, 1), (630, 1), (631, 1), (632, 1), (633, 1), (634,
1), (635, 1), (636, 1), (637, 1), (638, 1), (639, 1), (640, 1), (641,
1), (642, 1), (643, 1), (644, 1), (645, 1), (646, 1), (647, 1), (648, 1),
(649, 1), (650, 1), (651, 1), (652, 1), (653, 1), (654, 1), (655, 1), (656,
1), (657, 1), (658, 1), (659, 1), (660, 1), (661, 1), (662, 1), (663, 1),
(664, 1), (665, 1), (666, 1), (667, 1), (668, 1), (669, 1), (670, 1), (671,
1), (672, 1), (673, 1), (674, 1), (675, 1), (676, 1), (677, 1), (678, 1),
(679, 1), (680, 1), (681, 1), (682, 1), (683, 1), (684, 1), (685, 1), (686,
1), (687, 1), (688, 1), (689, 1), (690, 1), (691, 1), (692, 1), (693, 1),
(694, 1), (695, 1), (696, 1), (697, 1), (698, 1), (699, 1), (700, 1), (701,
1), (702, 1), (703, 1), (704, 1), (705, 1), (706, 1), (707, 1), (708, 1),
(709, 1), (710, 1), (711, 1), (712, 1), (713, 1), (714, 1), (715, 1), (716,
1), (717, 1), (718, 1), (719, 1), (720, 1), (721, 1), (722, 1), (723, 1),
(724, 1), (725, 1), (726, 1), (727, 1), (728, 1), (729, 1), (730, 1), (731,
1), (732, 1), (733, 1), (734, 1), (735, 1), (736, 1), (737, 1), (738, 1),
(739, 1), (740, 1), (741, 1), (742, 1), (743, 1), (744, 1), (745, 1), (746,
1), (747, 1), (748, 1), (749, 1), (750, 1), (751, 1), (752, 1), (753, 1),
(754, 1), (755, 1), (756, 1), (757, 1), (758, 1), (759, 1), (760, 1), (761,
1), (762, 1), (763, 1), (764, 1), (765, 1), (766, 1), (767, 1), (768, 1),
(769, 1), (770, 1), (771, 1), (772, 1), (773, 1), (774, 1), (775, 1), (776,
1), (777, 1), (778, 1), (779, 1), (780, 1), (781, 1), (782, 1), (783, 1),
(784, 1), (785, 1), (786, 1), (787, 1), (788, 1), (789, 1), (790, 1), (791,
1), (792, 1), (793, 1), (794, 1), (795, 1), (796, 1), (797, 1), (798, 1),
(799, 1), (800, 1), (801, 1), (802, 1), (803, 1), (804, 1), (805, 1), (806,
1), (807, 1), (808, 1), (809, 1), (810, 1), (811, 1), (812, 1), (813, 1),
(814, 1), (815, 1), (816, 1), (817, 1), (818, 1), (819, 1), (820, 1), (821,
1), (822, 1), (823, 1), (824, 1), (825, 1), (826, 1), (827, 1), (828, 1),
(829, 1), (830, 1), (831, 1), (832, 1), (833, 1), (834, 1), (835, 1), (836,
1), (837, 1), (838, 1), (839, 1), (840, 1), (841, 1), (842, 1), (843, 1),
(844, 1), (845, 1), (846, 1), (847, 1), (848, 1), (849, 1), (850, 1), (851,
1), (852, 1), (853, 1), (854, 1), (855, 1), (856, 1), (857, 1), (858, 1),
(859, 1), (860, 1), (861, 1), (862, 1), (863, 1), (864, 1), (865, 1), (866,
1), (867, 1), (868, 1), (869, 1), (870, 1), (871, 1), (872, 1), (873, 1),
(874, 1), (875, 1), (876, 1), (877, 1), (878, 1), (879, 1), (880, 1), (881,
1), (882, 1), (883, 1), (884, 1), (885, 1), (886, 1), (887, 1), (888, 1),
(889, 1), (890, 1), (891, 1), (892, 1), (893, 1), (894, 1), (895, 1), (896,
1), (897, 1), (898, 1), (899, 1), (900, 1), (901, 1), (902, 1), (903, 1),
(904, 1), (905, 1), (906, 1), (907, 1), (908, 1), (909, 1), (910, 1), (911,
1), (912, 1), (913, 1), (914, 1), (915, 1), (916, 1), (917, 1), (918, 1),
(919, 1), (920, 1), (921, 1), (922, 1), (923, 1), (924, 1), (925, 1), (926,
1), (927, 1), (928, 1), (929, 1), (930, 1), (931, 1), (932, 1), (933, 1),
(934, 1), (935, 1), (936, 1), (937, 1), (938, 1), (939, 1), (940, 1), (941,
1), (942, 1), (943, 1), (944, 1), (945, 1), (946, 1), (947, 1), (948, 1),
(949, 1), (950, 1), (951, 1), (952, 1), (953, 1), (954, 1), (955, 1), (956,
1), (957, 1), (958, 1), (959, 1), (960, 1), (961, 1), (962, 1), (963, 1),
(964, 1), (965, 1), (966, 1), (967, 1), (968, 1), (969, 1), (970, 1), (971,
1), (972, 1), (973, 1), (974, 1), (975, 1), (976, 1), (977, 1), (978, 1),
(979, 1), (980, 1), (981, 1), (982, 1), (983, 1), (984, 1), (985, 1), (986,
1), (987, 1), (988, 1), (989, 1), (990, 1), (991, 1), (992, 1), (993, 1),
(994, 1), (995, 1), (996, 1), (997, 1), (998, 1), (999, 1), (1000, 1)]
```

```
In [23]: import pandas as pd
import re
from collections import defaultdict

# Define a regular expression to tokenize the text
tokenizer = re.compile(r"\w+")

# Create an empty dictionary to store the inverted index
inverted_index = defaultdict(list)

# Loop over each document in the dataset
for i, row in df.iterrows():
    # Tokenize the text
    tokens = tokenizer.findall(row['descr'])
    # Count the frequency of each token
    term_freq = defaultdict(int)
    for token in tokens:
        term_freq[token] += 1
    # Add the document to the inverted index for each unique token
    for token, freq in term_freq.items():
        inverted_index[token].append((i, freq))

# Define a function to score documents based on a query
def score_documents(query, inverted_index, df):
    # Tokenize the query
    tokens = tokenizer.findall(query)
    # Count the frequency of each token
    query_freq = defaultdict(int)
    for token in tokens:
        query_freq[token] += 1
    # Compute the score for each document
    scores = defaultdict(int)
    for token, freq in query_freq.items():
        if token in inverted_index:
            for doc_id, doc_freq in inverted_index[token]:
                scores[doc_id] += freq * doc_freq
    # Sort the documents by score
    sorted_docs = sorted(scores.items(), key=lambda x: x[1], reverse=True)
    # Return the top-scoring documents as a DataFrame
    return df.iloc[[doc_id for doc_id, _ in sorted_docs]]

# Example usage: retrieve documents containing the word "election" and rank them
query = "UK election"
result_df = score_documents(query, inverted_index, df)
print(result_df.head())
```

Unnamed: 0 descr \

199	199	kennedy looks to election gains they may not ...
963	963	kennedys cautious optimism charles kennedy is...
1003	1003	labour mps fears over squabbling if there is ...
914	914	what the election should really be about a ge...
1000	1000	february poll claim speculation reports that ...

tags

199	politics, iraq, westminster hq, liberal democr...
963	politics, united kingdom, labour government, c...
1003	politics, mr brown, africa, brown camp, labour...
914	social issues, politics, mps, united kingdom, ...
1000	politics, baghdad, the sunday times, sunday te...

```
In [24]: # Define a regular expression to tokenize the text
tokenizer = re.compile(r"\w+")

# Create an empty dictionary to store the inverted index
inverted_index = defaultdict(list)

# Loop over each document in the dataset
for i, row in df.iterrows():
    # Tokenize the text
    tokens = tokenizer.findall(row['descr'])
    # Count the frequency of each token
    term_freq = defaultdict(int)
    for token in tokens:
        term_freq[token] += 1
    # Add the document to the inverted index for each unique token
    for token, freq in term_freq.items():
        inverted_index[token].append((i, freq))

# Define a function to score documents based on a query
def score_documents(query, inverted_index, df):
    # Tokenize the query
    tokens = tokenizer.findall(query)
    # Count the frequency of each token
    query_freq = defaultdict(int)
    for token in tokens:
        query_freq[token] += 1
    # Compute the score for each document
    scores = defaultdict(int)
    for token, freq in query_freq.items():
        if token in inverted_index:
            for doc_id, doc_freq in inverted_index[token]:
                scores[doc_id] += freq * doc_freq
    # Sort the documents by score
    sorted_docs = sorted(scores.items(), key=lambda x: x[1], reverse=True)
    # Return the top-scoring documents as a DataFrame
    return df.iloc[[doc_id for doc_id, _ in sorted_docs]]

# Example usage: retrieve documents containing the word "white" and ranking them
query = "white house"
result_df = score_documents(query, inverted_index, df)
print(result_df.head())
```

	Unnamed: 0	descr	tags
1836	1836	house prices rebound says halifax uk house pr...	business, halifax, london, halifax, capital ec...
1430	1430	white admits to balco drugs link banned ameri...	san francisco chronicle, la times, oil, remy k...
1629	1629	uk house prices dip in november uk house pric...	business, halifax, london, united kingdom, nor...
367	367	white prepared for battle toughscrummaging pr...	sports, zurich, cardiff, leicester, bristol, b...
859	859	no more concessions on terror charles clarke ...	politics, mps, bbc radio, bbc news, law lords,...

In the above result generated, we see that the documents retrieved are not very accurate for the query we specified as the results contain the token "white" or "house" and not both of them together. To solve this issue we use phrase queries, which we will see in the next sections.

```
In [25]: # Define a regular expression to tokenize the text
tokenizer = re.compile(r"\w+")

# Create an empty dictionary to store the inverted index
inverted_index = defaultdict(list)

# Loop over each document in the dataset
for i, row in df.iterrows():
    # Tokenize the text
    tokens = tokenizer.findall(row['descr'])
    # Count the frequency of each token
    term_freq = defaultdict(int)
    for token in tokens:
        term_freq[token] += 1
    # Add the document to the inverted index for each unique token
    for token, freq in term_freq.items():
        inverted_index[token].append((i, freq))

# Define a function to score documents based on a query
def score_documents(query, inverted_index, df):
    # Tokenize the query
    tokens = tokenizer.findall(query.lower())
    # Count the frequency of each token
    query_freq = defaultdict(int)
    for token in tokens:
        query_freq[token] += 1
    # Compute the score for each document
    scores = defaultdict(int)
    for token, freq in query_freq.items():
        if token in inverted_index:
            for doc_id, doc_freq in inverted_index[token]:
                scores[doc_id] += freq * doc_freq
    # Sort the documents by score
    sorted_docs = sorted(scores.items(), key=lambda x: x[1], reverse=True)
    # Return the top-scoring documents as a DataFrame
    return df.iloc[[doc_id for doc_id, _ in sorted_docs]]

# Example usage: retrieve documents containing the words "chelsea" and rank them
query = "chelsea"
result_df = score_documents(query, inverted_index, df)
print(result_df.head())
```

	Unnamed: 0	descr \
0	0	chelsea sack mutu chelsea have sacked adrian ...
12	12	desailly backs blues revenge trip marcel desa...
139	139	desailly backs blues revenge trip marcel desa...
155	155	chelsea hold arsenal a gripping game between ...
358	358	chelsea clinch cup in extratime after extrati...
		tags
0	0	sports, stamford bridge, football association,...
12	12	sports, barcelona, milan, the chelsea, bbc, eu...
139	139	sports, barcelona, milan, the chelsea, bbc, eu...
155	155	sports, henry, thierry henry, robert pires, wi...
358	358	sports, barcelona, liverpool, newcastle, reds,...

The above query is accurate as it retrieves all documents containing the word "chelsea" in it

```
In [26]: # Define a regular expression to tokenize the text
tokenizer = re.compile(r"\w+")

# Create an empty dictionary to store the inverted index
inverted_index = defaultdict(list)

# Loop over each document in the dataset
for i, row in df.iterrows():
    # Tokenize the text
    tokens = tokenizer.findall(row['descr'])
    # Count the frequency of each token
    term_freq = defaultdict(int)
    for token in tokens:
        term_freq[token] += 1
    # Add the document to the inverted index for each unique token
    for token, freq in term_freq.items():
        inverted_index[token].append((i, freq))

# Define a function to score documents based on a query
def score_documents(query, inverted_index, df):
    # Tokenize the query
    tokens = tokenizer.findall(query.lower())
    # Count the frequency of each token
    query_freq = defaultdict(int)
    for token in tokens:
        query_freq[token] += 1
    # Compute the score for each document
    scores = defaultdict(int)
    for token, freq in query_freq.items():
        if token in inverted_index:
            for doc_id, doc_freq in inverted_index[token]:
                scores[doc_id] += freq * doc_freq
    # Sort the documents by score
    sorted_docs = sorted(scores.items(), key=lambda x: x[1], reverse=True)
    # Return the top-scoring documents as a DataFrame
    return df.iloc[[doc_id for doc_id, _ in sorted_docs]]

# Example usage: retrieve documents containing the words "robot" and rank them
query = "robot"
result_df = score_documents(query, inverted_index, df)
print(result_df.head())
```

	Unnamed: 0	descr	tags
711	711	robots learn robotiquette rules robots are le...	technology, london, bbc, london's science muse...
426	426	humanoid robot learns how to run carmaker hon...	sony, honda, czech republic, electronics, car ...
780	780	hitachi unveils fastest robot japanese electr...	technology, washington dc, toyota, asimo, sony...
1244	1244	humanoid robot learns how to run carmaker hon...	spokesman, kelly holmes, king, europe, epo, ru...
721	721	gadget show heralds mp christmas partners of ...	entertainment, technology, london, sony, creat...

The above query is accurate as it retrieves all documents containing the word "robot" in it

```
In [27]: # Define a regular expression to tokenize the text
tokenizer = re.compile(r"\w+")

# Create an empty dictionary to store the inverted index
inverted_index = defaultdict(list)

# Loop over each document in the dataset
for i, row in df.iterrows():
    # Tokenize the text
    tokens = tokenizer.findall(row['descr'])
    # Count the frequency of each token
    term_freq = defaultdict(int)
    for token in tokens:
        term_freq[token] += 1
    # Add the document to the inverted index for each unique token
    for token, freq in term_freq.items():
        inverted_index[token].append((i, freq))

# Define a function to score documents based on a query
def score_documents(query, inverted_index, df):
    # Tokenize the query
    tokens = tokenizer.findall(query.lower())
    # Count the frequency of each token
    query_freq = defaultdict(int)
    for token in tokens:
        query_freq[token] += 1
    # Compute the score for each document
    scores = defaultdict(int)
    for token, freq in query_freq.items():
        if token in inverted_index:
            for doc_id, doc_freq in inverted_index[token]:
                scores[doc_id] += freq * doc_freq
    # Sort the documents by score
    sorted_docs = sorted(scores.items(), key=lambda x: x[1], reverse=True)
    # Return the top-scoring documents as a DataFrame
    return df.iloc[[doc_id for doc_id, _ in sorted_docs]]

# Example usage: retrieve documents containing the words "hollywood" and rank
query = "hollywood"
result_df = score_documents(query, inverted_index, df)
print(result_df.head())
```

	Unnamed: 0	descr	tags
1091	1091 godzilla gets hollywood fame star movie monst...		
2199	2199 keanu reeves given hollywood star actor keanu...		
754	754 games win for bluray dvd format the nextgener...		
762	762 games win for bluray dvd format the nextgener...		
469	469 movie body hits peertopeer nets the movie ind...		
			tags
1091	entertainment, human interest, walk of fame, g...		
2199	entertainment, human interest, beirut, volvo, ...		
754	entertainment, technology, las vegas, dell, to...		
762	entertainment, technology, las vegas, dell, to...		
469	entertainment, technology, law, phoenix, bitto...		

The above query is accurate as it retrieves all documents containing the word "hollywood" in it

```
In [28]: # Define a regular expression to tokenize the text
tokenizer = re.compile(r"\w+")

# Create an empty dictionary to store the inverted index
inverted_index = defaultdict(list)

# Loop over each document in the dataset
for i, row in df.iterrows():
    # Tokenize the text
    tokens = tokenizer.findall(row['descr'])
    # Count the frequency of each token
    term_freq = defaultdict(int)
    for token in tokens:
        term_freq[token] += 1
    # Add the document to the inverted index for each unique token
    for token, freq in term_freq.items():
        inverted_index[token].append((i, freq))

# Define a function to score documents based on a query
def score_documents(query, inverted_index, df):
    # Tokenize the query
    tokens = tokenizer.findall(query.lower())
    # Count the frequency of each token
    query_freq = defaultdict(int)
    for token in tokens:
        query_freq[token] += 1
    # Compute the score for each document
    scores = defaultdict(int)
    for token, freq in query_freq.items():
        if token in inverted_index:
            for doc_id, doc_freq in inverted_index[token]:
                scores[doc_id] += freq * doc_freq
    # Sort the documents by score
    sorted_docs = sorted(scores.items(), key=lambda x: x[1], reverse=True)
    # Return the top-scoring documents as a DataFrame
    return df.iloc[[doc_id for doc_id, _ in sorted_docs]]

# Example usage: retrieve documents containing the words "apple" and rank them
query = "appl"
result_df = score_documents(query, inverted_index, df)
print(result_df.head())
```

```
Empty DataFrame
Columns: [Unnamed: 0, descr, tags]
Index: []
```

We see that for the query "apple", we get the above results which are accurate. We can cross verify this by seeing the type of news this is in the tags fields, it says technology for the retrieved data.

```
In [29]: # Define a regular expression to tokenize the text
tokenizer = re.compile(r"\w+")

# Create an empty dictionary to store the inverted index
inverted_index = defaultdict(list)

# Loop over each document in the dataset
for i, row in df.iterrows():
    # Tokenize the text
    tokens = tokenizer.findall(row['descr'])
    # Count the frequency of each token
    term_freq = defaultdict(int)
    for token in tokens:
        term_freq[token] += 1
    # Add the document to the inverted index for each unique token
    for token, freq in term_freq.items():
        inverted_index[token].append((i, freq))

# Define a function to score documents based on a query
def score_documents(query, inverted_index, df):
    # Tokenize the query
    tokens = tokenizer.findall(query.lower())
    # Count the frequency of each token
    query_freq = defaultdict(int)
    for token in tokens:
        query_freq[token] += 1
    # Compute the score for each document
    scores = defaultdict(int)
    for token, freq in query_freq.items():
        if token in inverted_index:
            for doc_id, doc_freq in inverted_index[token]:
                scores[doc_id] += freq * doc_freq
    # Sort the documents by score
    sorted_docs = sorted(scores.items(), key=lambda x: x[1], reverse=True)
    # Return the top-scoring documents as a DataFrame
    return df.iloc[[doc_id for doc_id, _ in sorted_docs]]

# Example usage: retrieve documents containing the words "artificial" and "int
query = "blockchain"
result_df = score_documents(query, inverted_index, df)
print(result_df.head())
```

Empty DataFrame
Columns: [Unnamed: 0, descr, tags]
Index: []

For this particular query, we do not have any document in the dataset which has the word "blockchain". Hence the retrieved dataframe is empty output for the query and would be an empty result set. That is, the query would return no documents, since there are no documents that match the query.

```
In [30]: # Define a regular expression to tokenize the text
tokenizer = re.compile(r"\w+")

# Create an empty dictionary to store the inverted index
inverted_index = defaultdict(list)

# Loop over each document in the dataset
for i, row in df.iterrows():
    # Tokenize the text
    tokens = tokenizer.findall(row['descr'])
    # Count the frequency of each token
    term_freq = defaultdict(int)
    for token in tokens:
        term_freq[token] += 1
    # Add the document to the inverted index for each unique token
    for token, freq in term_freq.items():
        inverted_index[token].append((i, freq))

# Define a function to score documents based on a query
def score_documents(query, inverted_index, df):
    # Tokenize the query
    tokens = tokenizer.findall(query.lower())
    # Count the frequency of each token
    query_freq = defaultdict(int)
    for token in tokens:
        query_freq[token] += 1
    # Compute the score for each document
    scores = defaultdict(int)
    for token, freq in query_freq.items():
        if token in inverted_index:
            for doc_id, doc_freq in inverted_index[token]:
                scores[doc_id] += freq * doc_freq
    # Sort the documents by score
    sorted_docs = sorted(scores.items(), key=lambda x: x[1], reverse=True)
    # Return the top-scoring documents as a DataFrame
    return df.iloc[[doc_id for doc_id, _ in sorted_docs]]

# Example usage: retrieve documents containing the words "artificial" and "int
query = "pear"
result_df = score_documents(query, inverted_index, df)
print(result_df.head())
```

Empty DataFrame
Columns: [Unnamed: 0, descr, tags]
Index: []

When we searched for the query "apple" we got documents which matched with it. But when we tried it with a different fruit, "pear" we have no matching documents in the dataset

Section 4.3

Result of Wild Card queries

In this section, we have used the asterisks symbol at multiple places in the query. The description of each is given below the output of the code snippets

```
In [31]: # read in the dataset
df = pd.read_csv("/content/BBCNews.csv")

# define the wildcard search function
def wildcard_search(query, df):
    # replace the wildcard character with a regular expression
    query = query.replace("*", ".*")
    query = query.replace("?", ".") 

    # search the text column of the dataframe using the regular expression
    results = df[df["descr"].str.contains(query, case=False)].index

    return results

# test the wildcard search function
query = "robo*"
results = wildcard_search(query, df)
print(results)
```

```
Int64Index([ 124,  426,  492,  500,  541,  544,  579,  588,  711,  721,  728,
             734,  780,  920,  1025, 1232, 1244, 1310, 1318, 1359, 1362, 1397,
             1406, 2071, 2098, 2118, 2250],
            dtype='int64')
```

The above function will return a list of document indices that contain words starting with "robo", such as "robot" and "robotics".

In [32]:

```
# read in the dataset
df = pd.read_csv("/content/BBCNews.csv")

# define the wildcard search function
def wildcard_search(query, df):
    # replace the wildcard character with a regular expression
    query = query.replace("*", ".*")
    query = query.replace("?", ".") 

    # search the text column of the dataframe using the regular expression
    results = df[df["descr"].str.contains(query, case=False)].index

    return results

# test the wildcard search function
query = "*tech"
results = wildcard_search(query, df)
print(results)
```

```
Int64Index([ 37,   38,   84,   92,   95,  117,  152,  204,  238,  256,
             ...
            2293, 2299, 2317, 2320, 2339, 2349, 2358, 2366, 2369, 2398],
           dtype='int64', length=420)
```

The above query will match all documents that end with the word "tech". This could be useful for finding articles about different types of technology or tech-related news.

In [33]:

```
# read in the dataset
df = pd.read_csv("/content/BBCNews.csv")

# define the wildcard search function
def wildcard_search(query, df):
    # replace the wildcard character with a regular expression
    query = query.replace("*", ".*")
    query = query.replace("?", ".") 

    # search the text column of the dataframe using the regular expression
    results = df[df["descr"].str.contains(query, case=False)].index

    return results

# test the wildcard search function
query = "cyber*attack"
results = wildcard_search(query, df)
print(results)
```

```
Int64Index([474, 482, 516, 525, 554, 786, 1292, 1300, 1334, 1343, 1372], dtype='int64')
```

The above query will match all documents that contain words starting with "cyber" and ending with "attack". This could be useful for finding articles about cybersecurity, cybercrime, or cyber attacks.

```
In [34]: # read in the dataset
df = pd.read_csv("/content/BBCNews.csv")

# define the wildcard search function
def wildcard_search(query, df):
    # replace the wildcard character with a regular expression
    query = query.replace("*", ".*")
    query = query.replace("?", ".") 

    # search the text column of the dataframe using the regular expression
    results = df[df["descr"].str.contains(query, case=False)].index

    return results
```

```
# test the wildcard search function
query = "mark*crash"
results = wildcard_search(query, df)
print(results)
```

```
Int64Index([22, 92, 152, 181, 495, 573, 807, 822, 1313, 1391, 1516, 1580, 179
0,
           1969],
           dtype='int64')
```

```
In [35]: # read in the dataset
df = pd.read_csv("/content/BBCNews.csv")

# define the wildcard search function
def wildcard_search(query, df):
    # replace the wildcard character with a regular expression
    query = query.replace("*", ".*")
    query = query.replace("?", ".") 

    # search the text column of the dataframe using the regular expression
    results = df[df["descr"].str.contains(query, case=False)].index

    return results

# test the wildcard search function
query = "*energy*"
results = wildcard_search(query, df)
print(results)
```

```
Int64Index([ 21,  57, 161, 162, 191, 212, 221, 238, 253, 308, 345,
 351, 359, 360, 420, 487, 571, 572, 689, 723, 741, 746,
 783, 813, 862, 973, 1216, 1305, 1389, 1390, 1433, 1434, 1439,
 1472, 1516, 1575, 1607, 1634, 1678, 1685, 1695, 1705, 1712, 1731,
 1738, 1745, 1770, 1788, 1792, 1812, 1837, 1847, 1853, 1867, 1876,
 1895, 1910, 1914, 1927, 1929, 1931, 1934, 2007, 2211, 2227, 2232,
 2233, 2235, 2244, 2278, 2291, 2312, 2343, 2362, 2372, 2389, 240
8],
           dtype='int64')
```

The above query will match all documents that contain the word "energy" or any word that starts or ends with "energy". This could be useful for finding articles about renewable energy, fossil fuels, or energy policy.

Section 4.4

Result of Phrase queries:

In free text queries we saw that if either of the words were present in the query, we got a result. In phrase queries, unlike free txt queries the query given has to be in the same format to find a match.

```
In [36]: # Define a phrase to search for
phrase = 'climate change'

# Find all documents that contain the phrase
matches = df[df['descr'].str.contains(phrase)]

# Print the number of documents that match the query
print(f"Found {len(matches)} documents containing the phrase '{phrase}':")

# Loop over the matching documents and print their IDs, titles, and text
for i, row in matches.iterrows():
    print(f"\nDocument ID: {i}")
    print(f"Article: {row['descr']}")
```

```
Found 11 documents containing the phrase 'climate change':
```

```
Document ID: 221
```

```
Article: uk set to cut back on embassies nine overseas embassies and high
commissions will close in an effort to save money uk foreign secretary jac
k straw has announced the bahamas east timor madagascar and swaziland are
among the areas affected by the biggest shakeup for the diplomatic service
for years other diplomatic posts are being turned over to local staff mr s
traw said the move would save m a year to free up cash for priorities such
as fighting terrorism honorary consuls will be appointed in some of the a
reas affected by the embassy closures nine consulates or consulates genera
l will also be closed mostly in europe and america they include dallas in
the us bordeaux in france and oporto in portugal with local staff replacin
g uk representation in another the changes are due to be put in place bef
ore the end of with most savings made from cutting staff and running cost
s some of the money will have to be used to fund redundancy payments in a
written statement mr straw said the savings made will help to underpin hig
her priority work in line with the foreign and commonwealth offices strate
gic priorities including counter proliferation counterterrorism energy and
```

```
In [37]: # Define a phrase to search for
phrase = 'gender pay gap'

# Find all documents that contain the phrase
matches = df[df['descr'].str.contains(phrase)]

# Print the number of documents that match the query
print(f"Found {len(matches)} documents containing the phrase '{phrase}':")

# Loop over the matching documents and print their IDs, titles, and text
for i, row in matches.iterrows():
    print(f"\nDocument ID: {i}")
    print(f"Article: {row['descr']}")
```

Found 1 documents containing the phrase 'gender pay gap':

Document ID: 1161

Article: hewitt decries career sexism plans to extend paid maternity leave beyond six months should be prominent in labours election manifesto the trade and industry secretary has said patricia hewitt said the cost of the proposals was being evaluated but it was an increasingly high priority and a shared goal across government ms hewitt was speaking at a gender and productivity seminar organised by the equal opportunities commission eoc mothers can currently take up to six months paid leave and six unpaid ms hewitt told the seminar clearly one of the things we need to do in the future is to extend the period of payment for maternity leave beyond the first six months into the second six months we are looking at how quickly we can do that because obviously there are cost implications because the taxpayer reimburses the employers for the cost of that ms hewitt also announced a new drive to help women who want to work in male dominated sectors saying sexism at work was still preventing women reaching their full potential plans include funding for universities to help female science and engineering graduates find jobs and taster courses for men and women in nontraditional jobs women in fulltime work earn less than men according to the equal opportunities commission eoc the minister told delegates that getting rid of career sexism was vital to closing the gender pay gap career sexism limits opportunities for women of all ages and prevents them from achieving their full potential it is simply wrong to assume someone cannot do a job on the grounds of their sex she said earlier she told bbc radio 4's today programme what we are talking about here is the fact that about six out of women work in jobs that are lowpaid and typically dominated by women so we have got very segregated employment unfortunately in some cases this reflects very oldfashioned and stereotypical ideas about the appropriate jobs for women or indeed for men career sexism is about saying that engineering for instance where only of employees are women is really a maledominated industry construction is even worse but it is also about saying childcare jobs are really there for women and not suitable for men career sexism goes both ways she added that while progress had been made there was still a gap in pay figures the average woman working fulltime is being paid about p for every pound a man is earning for women working parttime it is p the department for trade and industry will also provide funding to help a new pay experts panel run by the tuc it has been set up to advise hundreds of companies on equal wage policies research conducted by the eoc last year revealed that many britons believe the pay gap between men and women is the result of natural differences between the sexes women hold less than of the top positions in ftse companies the police the judiciary and trade unions according to their figures and retired women have just over half the income of their male counterparts on average

Tags: labor, bbc radio, ftse 100, equal opportunities commission, department for trade and industry, patricia hewitt, industry secretary, minister, bbc radio 4, gender equality, family law, discrimination, sexism, family, politics, chauvinism, human behavior, patricia hewitt, parental leave, maternity leave in the united states, gender pay gap in russia

```
In [38]: # Define a phrase to search for
phrase = 'global economic'

# Find all documents that contain the phrase
matches = df[df['descr'].str.contains(phrase)]

# Print the number of documents that match the query
print(f"Found {len(matches)} documents containing the phrase '{phrase}':")

# Loop over the matching documents and print their IDs, titles, and text
for i, row in matches.iterrows():
    print(f"\nDocument ID: {i}")
    print(f"Article: {row['descr']}")
```

Found 7 documents containing the phrase 'global economic':

Document ID: 1690

Article: newest eu members underpin growth the european unions newest members will bolster europe's economic growth in according to a new report the eight central european states which joined the eu last year will see growth the united nations economic commission for europe unece said in contrast the euro zone countries will put in a lacklustre performance generating growth of only the global economy will slow in the unece forecasts due to widespread weakness in consumer demand it warned that growth could also be threatened by attempts to reduce the united states huge current account deficit which in turn might lead to significant volatility in exchange rates unece is forecasting average economic growth of across the european union in however total output across the euro zone is forecast to fall in from to this is due largely to the faltering german economy which shrank in the last quarter of on monday germanys bdb private banks association said the german economy would struggle to meet its growth target in separately the bundesbank warned that germanys efforts to reduce its budget deficit below of gdp presented huge risks given that headline econo

```
In [39]: # Define a phrase to search for
phrase = 'food waste reduction initiatives'

# Find all documents that contain the phrase
matches = df[df['descr'].str.contains(phrase)]

# Print the number of documents that match the query
print(f"Found {len(matches)} documents containing the phrase '{phrase}':")

# Loop over the matching documents and print their IDs, titles, and text
for i, row in matches.iterrows():
    print(f"\nDocument ID: {i}")
    print(f"Article: {row['descr']}")
```

Found 0 documents containing the phrase 'food waste reduction initiatives':

The above phrase has no documents found as no document in the dataset has the above phrase in it

```
In [40]: # Define a phrase to search for
phrase = 'human rights'

# Find all documents that contain the phrase
matches = df[df['descr'].str.contains(phrase)]

# Print the number of documents that match the query
print(f"Found {len(matches)} documents containing the phrase '{phrase}'")

# Loop over the matching documents and print their IDs, titles, and text
for i, row in matches.iterrows():
    print(f"\nDocument ID: {i}")
    print(f"Article: {row['descr']}")
```

Found 55 documents containing the phrase 'human rights':

```
Document ID: 216
Article: custody death rate shocks mps deaths in custody have reached shocking levels a committee of mps and peers has warned the joint committee on human rights found those committing suicide were mainly the most vulnerable with mental health drugs or alcohol problems members urged the government to set up a task force to tackle deaths in prisons police cells detention centres and special hospitals there was one prison suicide every four days between and mps said the report which followed a yearlong inquiry by the committee found the high death rate amounts to a serious failure to protect the right to life of a highly vulnerable group many of those who ended up taking their own lives had presented themselves to the authorities with these problems before they even offended the report said it questioned whether prison was the most appropriate place for them to be kept and whether earlier intervention would have meant custody could have been avoided increased resources and a reduction in the use of imprisonment was needed to address the issue in the longer term the report said committee chairman labour mp jean corston said each and every death in custody is a dea
```

The above output is very large. We have used the below code to retrieve the 10 most relevant documents containing "human rights" in them. We have used TF-IDF and cosine similarity to get this output

```
In [41]: # Load the dataset into a pandas dataframe
df = pd.read_csv("/content/BBCNews.csv")

# Create a document-term matrix using TfidfVectorizer
vectorizer = TfidfVectorizer()
doc_term_matrix = vectorizer.fit_transform(df['descr'])

# Define the phrase query and perform the search
query = "human rights"
query_vec = vectorizer.transform([query])
doc_scores = cosine_similarity(query_vec, doc_term_matrix)[0]
indices = np.argsort(doc_scores)[::-1]

# Print the top 10 documents with matching phrase and their document ids
print("Top 10 documents with matching phrase:")
for i in range(10):
    doc_id = indices[i]
    doc_text = df.iloc[doc_id]['descr']
    doc_score = doc_scores[doc_id]
    if "human rights" in doc_text:
        print(f"Document ID: {doc_id}, Score: {doc_score}")
        print(f'Article: {doc_text}')
        print("-----")
```

Top 10 documents with matching phrase:

Document ID: 830, Score: 0.6321156770132687

Article: amnesty chief laments war failure the lack of public outrage about the war on terror is a powerful indictment of the failure of human rights groups amnesty internationals chief has said in a lecture at the london school of economics irene khan said human rights had been flouted in the name of security since september she said the human rights movement had to use simpler language both to prevent scepticism and spread a moral message and it had to fight poverty not just focus on political rights for elites ms khan highlighted detentions without trial including those at the us camp at guantanamo bay in cuba and the abuse of prisoners as evidence of increasing human rights problems whats a new challenge is the way in which this ageold debate on security and human rights has been translated into the language of war she said by using the language of war human rights are being sidelined because we know human rights do not apply in times of war ms khan said such breaches were infectious and were now seen in almost very major country in the world the human rights movement faces a crisis of faith in the value of human rights she said that was accompanied by a crisis of governance where the united nations system did not seem able to hold

```
In [42]: import pandas as pd
import numpy as np
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics.pairwise import cosine_similarity

# Load the dataset into a pandas dataframe
#df = pd.read_csv("BBC News Dataset.csv")

# Create a document-term matrix using TfidfVectorizer
vectorizer = TfidfVectorizer()
doc_term_matrix = vectorizer.fit_transform(df['descr'])

# Define the phrase query and perform the search
query = "economic growth"
query_vec = vectorizer.transform([query])
doc_scores = cosine_similarity(query_vec, doc_term_matrix)[0]
indices = np.argsort(doc_scores)[::-1]

# Print the top 10 documents with matching phrase and their document ids
print("Top 10 documents with matching phrase:")
for i in range(10):
    doc_id = indices[i]
    doc_text = df.iloc[doc_id]['descr']
    doc_score = doc_scores[doc_id]
    if "economic growth" in doc_text:
        print(f"Document ID: {doc_id}, Score: {doc_score}")
        print(f'Article: {doc_text}')
        print("-----")
```

Top 10 documents with matching phrase:

Document ID: 1690, Score: 0.4048869826479714

Article: newest eu members underpin growth the european unions newest members will bolster europe's economic growth in according to a new report the eight central european states which joined the eu last year will see growth the united nations economic commission for europe unece said in contrast the euro zone countries will put in a lacklustre performance generating growth of only the global economy will slow in the unece forecasts due to widespread weakness in consumer demand it warned that growth could also be threatened by attempts to reduce the united states huge current account deficit which in turn might lead to significant volatility in exchange rates unece is forecasting average economic growth of across the euro union in however total output across the euro zone is forecast to fall in from to this is due largely to the faltering german economy which shrank in the last quarter of on monday germanys bdb private banks association said the german economy would struggle to meet its growth target in separately the bundesbank warned that germanys efforts to reduce its budget deficit below of gdp presented huge risks given that headline economic growth was set to fall below this year publishing its economic survey

The code defines the phrase query "climate change action" and performs the search using cosine similarity between the query vector and the document-term matrix. The code prints the top 10 documents with matching phrase and their document IDs, ranked by the cosine similarity score.

Section 4.5

How the evaluator can test, an arbitrary text query relevant to your corpus?

Ans: For an evaluator to test an arbitrary text query, we can use any of the above mentioned methods, free text, wild card or phrase queries. One method to do this is by taking an input from the user and performing the above code snippets from section 4.1, 4.2, 4.3 and 4.4. The other way is we manually plug in the input query given by the user and run the code. Both ways, we can figure out whether the query is relevant to our corpus or not. It acts as a search engine for a specified dataset. If the query is relevant, the output is a set of documents containing the query. We can further enhance the relevance of the query by using cosine similarity, TF-IDF, ranking and other parameters to retrieve the most relevant documents. If the output is empty, we know that the query provided is irrelevant to the corpus.

In greater detail, we can follow the below steps to evaluate an arbitrary text query.

To evaluate the relevance of an arbitrary text query to the corpus, you can use a common evaluation metric in information retrieval called "precision at k" ($P@k$).

Here's how it works:

- Choose a value for k , which represents the number of documents to consider for each query.
- For each query, retrieve the top k documents from the corpus that match the query. You can use any information retrieval technique such as TF-IDF, BM25, or neural models for this.
- Evaluate the relevance of each retrieved document using some relevance criteria such as relevance judgments made by human assessors.
- Compute the precision at k by dividing the number of relevant documents retrieved by the total number of documents retrieved up to k .
- Repeat the above steps for multiple queries to obtain an average $P@k$ score.

By using precision at k , you can evaluate how well a retrieval system is performing for a given query. The higher the $P@k$ score, the more relevant the retrieved documents are for the query.

The choice of k depends on the specific application and the available resources. A larger k would give more comprehensive results but would also increase the computational cost.

Section 4.6

Any one additional functionality: Relevance feedback, Semantic matching, re-ranking of results, and finding out query intention.

Semantic matching is a technique that involves identifying the underlying meaning of the query and matching it with the meaning of the documents in the corpus. This can be done using various natural language processing (NLP) techniques such as semantic analysis, named entity recognition, and dependency parsing.

The steps below can be followed to check:

- Load the dataset: You can load the BBC News dataset into a pandas dataframe using the following code:

```
In [43]: import pandas as pd  
  
data = pd.read_csv('/content/BBCNews.csv')
```

- Preprocess the data: Before performing any retrieval or semantic matching, you need to preprocess the data. This involves converting the text data into a format that can be used for retrieval, such as tokenizing the text, removing stop words, and stemming the words.

```
In [44]: import nltk  
from nltk.corpus import stopwords  
from nltk.tokenize import word_tokenize  
from nltk.stem.porter import PorterStemmer  
  
nltk.download('stopwords')  
nltk.download('punkt')  
  
stemmer = PorterStemmer()  
stop_words = set(stopwords.words('english'))  
  
def preprocess_text(text):  
    tokens = word_tokenize(text.lower())  
    tokens = [stemmer.stem(token) for token in tokens if token not in stop_words]  
    return ' '.join(tokens)  
  
data['processed_text'] = data['descr'].apply(preprocess_text)
```



```
[nltk_data] Downloading package stopwords to /root/nltk_data...  
[nltk_data]   Package stopwords is already up-to-date!  
[nltk_data] Downloading package punkt to /root/nltk_data...  
[nltk_data]   Package punkt is already up-to-date!
```

- Extract semantic features: To perform semantic matching, you need to extract semantic features from the text data. One way to do this is to use pre-trained word embeddings, such as Word2Vec or GloVe.

```
In [45]: import gensim.downloader as api
import numpy as np

word_vectors = api.load("glove-wiki-gigaword-100")

def extract_semantic_features(text):
    tokens = word_tokenize(text.lower())
    vectors = []
    for token in tokens:
        if token in word_vectors:
            vectors.append(word_vectors[token])
    if vectors:
        return np.mean(vectors, axis=0)
    else:
        return np.zeros(100)

data['semantic_features'] = data['processed_text'].apply(extract_semantic_feat
```

- Retrieve documents using a query: You can use semantic matching to retrieve documents based on a query by computing the cosine similarity between the semantic features of the query and the semantic features of the documents. For example, you can use the following code to retrieve the top 10 documents related to the query "technology":

```
In [46]: import numpy as np
```

```
query = 'technology'
query_features = extract_semantic_features(query)

scores = data['semantic_features'].apply(lambda x: np.dot(x, query_features))
ranked_documents = sorted(enumerate(scores), key=lambda x: x[1], reverse=True)

for document_index, score in ranked_documents:
    print(data.iloc[document_index]['descr'])
```

electronics firms eye plasma deal consumer electronics giants hitachi and matshushita electric are joining forces to share and develop technology for flat screen televisions the tieup comes as the worlds top producers are having to contend with falling prices and intense competition the two japanese companies will collaborate in research development production marketing and licensing they said the agreement would enable the two companies to expand the plasma display tv market globally plasma display panels are used for large thin tvs which are replacing oldstyle televisions the display market for highdefinition televisions is split between models using plasma display panels and others manufactured by the likes of sony and samsung using liquidcrystal displays lcds the deal will enable hitachi and matsushita which makes panasonic brand products to develop new technology and improve their competitiveness hitachi recently announced a deal to buy plasma display technology from rival fujitsu in an effort to strengthen its presence in the market separately fujitsu announced on monday that it is quitting the lcd panel market by transferring its operations in the area to japanese manufacturer sharp sharp will inherit staff manufacturing facilities and intellectual property from fujitsu the plasma panel market has seen rapid consolidation in recent months as the price of consumer electronic

- Improve the performance: You can improve the performance of the retrieval system by using various techniques such as relevance feedback or query expansion. You can also try different pre-trained word embeddings or fine-tune your own word embeddings using the dataset.