

```
nltk: natural language toolkit
```

```
In [1]: 1 # !pip install nltk
```

```
In [2]: 1 import nltk
```

## Important libraries

```
In [3]: 1 #for tokenization
2 from nltk import word_tokenize, sent_tokenize, WhitespaceTokenizer
3 #for Lemmitization
4 from nltk import WordNetLemmatizer, LancasterStemmer
5 #for handling stopwords
6 from nltk.corpus import stopwords
7 #frequency distribution
8 from nltk.probability import FreqDist
9 #Handling punctuations
10 from string import punctuation
11 #for contraction mapping
12 import contractions
13 #for wordcloud
14 from wordcloud import WordCloud
```

```
In [5]: 1 Text = """
2 Natural language processing (NLP) is a field of artificial intelligence (AI) that focuses on the int
3 computers and humans through natural language. The ultimate objective of NLP is to enable computers
4 interpret, and generate human language in a way that is both meaningful and useful.NLP techniques ar
5 text, extract relevant information, and derive insights from unstructured data.
6 """
```

```
In [6]: 1 #word tokenization
2 word_token=word_tokenize(Text)
3 # word_token
```

```
In [7]: 1 #sentence tokenization
2 sent_token=sent_tokenize(Text)
3 sent_token
```

```
Out[7]: ['\nNatural language processing (NLP) is a field of artificial intelligence (AI) that focuses on the in
teraction between \ncomputers and humans through natural language.',
 'The ultimate objective of NLP is to enable computers to understand, \ninterpret, and generate human l
anguage in a way that is both meaningful and useful.NLP techniques are used to analyze \ntext, extract
relevant information, and derive insights from unstructured data.']
```

```
In [8]: 1 #to tokenize on the basis of whitespace
        2 whitespace_token=WhitespaceTokenizer().tokenize(Text)
        3 whitespace_token
```

```
Out[8]: ['Natural',
         'language',
         'processing',
         '(NLP)',
         'is',
         'a',
         'field',
         'of',
         'artificial',
         'intelligence',
         '(AI)',
         'that',
         'focuses',
         'on',
         'the',
         'interaction',
         'between',
         'computers',
         'and',
         'humans',
         'through',
         'natural',
         'language.',
         'The',
         'ultimate',
         'objective',
         'of',
         'NLP',
         'is',
         'to',
         'enable',
         'computers',
         'to',
         'understand,',
         'interpret,',
         'and',
         'generate',
         'human',
         'language',
         'in',
         'a',
         'way',
         'that',
         'is',
         'both',
         'meaningful',
         'and',
         'useful.NLP',
         'techniques',
         'are',
         'used',
         'to',
         'analyze',
         'text,',
         'extract',
         'relevant',
         'information,',
         'and',
         'derive',
         'insights',
         'from',
         'unstructured',
         'data.']
```

```
In [9]: 1 #formatting text to same case
        2 lower_text=[x.lower() for x in word_token]
        3 lower_text
```

```
Out[9]: ['natural',
         'language',
         'processing',
         '(',
         'nlp',
         ')',
         'is',
         'a',
         'field',
         'of',
         'artificial',
         'intelligence',
         '(',
         'ai',
         ')',
         'that',
         'focuses',
         'on',
         'the',
         'interaction',
         'between',
         'computers',
         'and',
         'humans',
         'through',
         'natural',
         'language',
         ', ',
         'the',
         'ultimate',
         'objective',
         'of',
         'nlp',
         'is',
         'to',
         'enable',
         'computers',
         'to',
         'understand',
         ', ',
         'interpret',
         ', ',
         'and',
         'generate',
         'human',
         'language',
         'in',
         'a',
         'way',
         'that',
         'is',
         'both',
         'meaningful',
         'and',
         'useful.nlp',
         'techniques',
         'are',
         'used',
         'to',
         'analyze',
         'text',
         ', ',
         'extract',
         'relevant',
         'information',
         ', ',
         'and',
         'derive',
         'insights',
         'from',
         'unstructured',
```

```
'data',  
'.']
```

```
In [10]: 1 #Handling Punctuations
          2 punct_rem=punctuation
          3 print(punct_rem)
          4 punct_rem1=[word for word in lower_text if word not in punctuation ]
          5 punct_rem1
```

```
!"#$%&'()*+,-./:;<=>?@[\\]^_`{|}~
```

```
Out[10]: ['natural',
          'language',
          'processing',
          'nlp',
          'is',
          'a',
          'field',
          'of',
          'artificial',
          'intelligence',
          'ai',
          'that',
          'focuses',
          'on',
          'the',
          'interaction',
          'between',
          'computers',
          'and',
          'humans',
          'through',
          'natural',
          'language',
          'the',
          'ultimate',
          'objective',
          'of',
          'nlp',
          'is',
          'to',
          'enable',
          'computers',
          'to',
          'understand',
          'interpret',
          'and',
          'generate',
          'human',
          'language',
          'in',
          'a',
          'way',
          'that',
          'is',
          'both',
          'meaningful',
          'and',
          'useful.nlp',
          'techniques',
          'are',
          'used',
          'to',
          'analyze',
          'text',
          'extract',
          'relevant',
          'information',
          'and',
          'derive',
          'insights',
          'from',
          'unstructured',
          'data']
```

```
In [11]: 1 #Contraction mapping
2 text1= "We can't like to attend lecture"
3 contract_text=contractions.fix(text1)
4 contract_text
```

Out[11]: 'We cannot like to attend lecture'

```
In [12]: 1 stopwords_list=stopwords.words("English")
2 print(stopwords_list)
3 stopwords_rem1=[word for word in punct_rem1 if word not in stopwords_list ]
4 stopwords_rem1
```

```
['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", "you'll", "you
u'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', 'she', "she's", 'her',
'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them', 'their', 'theirs', 'themselves', 'wha
t', 'which', 'who', 'whom', 'this', 'that', "that'll", 'these', 'those', 'am', 'is', 'are', 'was', 'wer
e', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'do', 'does', 'did', 'doing', 'a', 'an', 'th
e', 'and', 'but', 'if', 'or', 'because', 'as', 'until', 'while', 'of', 'at', 'by', 'for', 'with', 'abou
t', 'against', 'between', 'into', 'through', 'during', 'before', 'after', 'above', 'below', 'to', 'fro
m', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'again', 'further', 'then', 'once', 'her
e', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', 'few', 'more', 'most', 'othe
r', 'some', 'such', 'no', 'nor', 'not', 'only', 'own', 'same', 'so', 'than', 'too', 'very', 's', 't',
'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now', 'd', 'll', 'm', 'o', 're', 've',
'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'doesn', "doesn't", 'hadn', "had
n't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn', "mightn't", 'mustn', "must
n't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn', "wasn't", 'weren', "were
n't", 'won', "won't", 'wouldn', "wouldn't"]
```

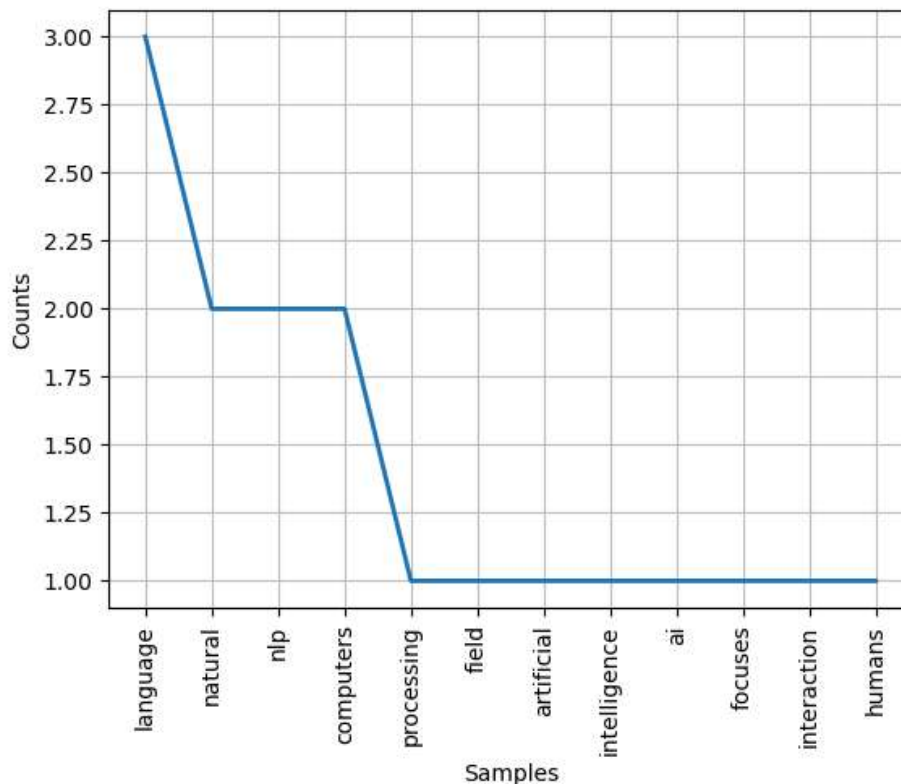
```
Out[12]: ['natural',
'language',
'processing',
'nlp',
'field',
'artificial',
'intelligence',
'ai',
'focuses',
'interaction',
'computers',
'humans',
'natural',
'language',
'ultimate',
'objective',
'nlp',
'enable',
'computers',
'understand',
'interpret',
'generate',
'human',
'language',
'way',
'meaningful',
'useful.nlp',
'techniques',
'used',
'analyze',
'text',
'extract',
'relevant',
'information',
'derive',
'insights',
'unstructured',
'data']
```

```
In [13]: 1 frequency_distribution=FreqDist(stopwords_rem1)
2 print(frequency_distribution)
3
4 # To find most common words using Frequency Distribution, add the following lines in above code :
5 print(frequency_distribution.most_common())
```

<FreqDist with 33 samples and 38 outcomes>

```
[('language', 3), ('natural', 2), ('nlp', 2), ('computers', 2), ('processing', 1), ('field', 1), ('arti-
ficial', 1), ('intelligence', 1), ('ai', 1), ('focuses', 1), ('interaction', 1), ('humans', 1), ('ultim-
ate', 1), ('objective', 1), ('enable', 1), ('understand', 1), ('interpret', 1), ('generate', 1), ('huma-
n', 1), ('way', 1), ('meaningful', 1), ('useful.nlp', 1), ('techniques', 1), ('used', 1), ('analyze',
1), ('text', 1), ('extract', 1), ('relevant', 1), ('information', 1), ('derive', 1), ('insights', 1),
('unstructured', 1), ('data', 1)]
```

```
In [14]: 1 import matplotlib.pyplot as plt
2 frequency_distribution.plot(12,cumulative=False)
3 plt.show()
```



```
In [15]: 1 text_lem=WordNetLemmatizer()
2 text_stem=LancasterStemmer()
3 for word in stopwords_rem1:
4     text_lem1=text_lem.lemmatize(word)
5     text_stem1=text_stem.stem(word)
6     print(f"original text: {word}")
7     print(f"lemmitize text: {text_lem1}")
8     print(f"steming text: {text_stem1}")
9     print(".*20")
```

```
original text: information
lemmitize text: information
steming text: inform
.....
original text: derive
lemmitize text: derive
steming text: der
.....
original text: insights
lemmitize text: insight
steming text: insight
.....
original text: unstructured
lemmitize text: unstructured
steming text: unstruct
.....
original text: data
lemmitize text: data
steming text: dat
.....
```

```
In [16]: 1 from nltk.sentiment.vader import SentimentIntensityAnalyzer
2 obj=SentimentIntensityAnalyzer()
3 Text="I like tea"
4 sent_text=obj.polarity_scores(Text)
5 sent_text
```

```
Out[16]: {'neg': 0.0, 'neu': 0.286, 'pos': 0.714, 'compound': 0.3612}
```

```
In [17]: 1 Text="very bad!"
2 sent_text=obj.polarity_scores(Text)
3 sent_text
```

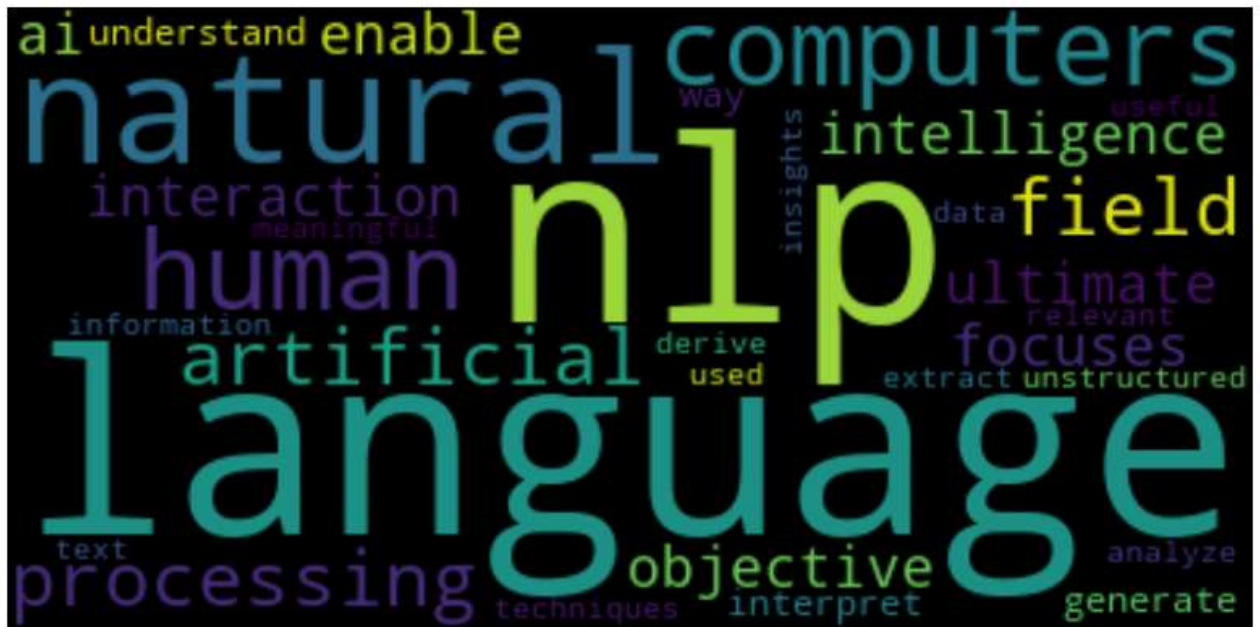
```
Out[17]: {'neg': 0.803, 'neu': 0.197, 'pos': 0.0, 'compound': -0.623}
```

```
In [18]: 1 from wordcloud import WordCloud
2 def plot_wordcloud(stopwords_rem1):
3     text = ' '.join(stopwords_rem1)
4     wordcloud = WordCloud(width=500, height=300, background_color='white').generate(text)
5     plt.figure(figsize=(10, 5))
6     plt.imshow(wordcloud, interpolation='bilinear')
7     # plt.axis('off')
8     plt.show()
```

## wordcloud

```
In [19]: 1
2 text = ' '.join(stopwords_rem1) #processed text(after removing stopwords,punctuations i.e. noise )
3 print(text)
4 wordcloud = WordCloud().generate(text)
5 plt.figure(figsize=(10, 5))
6 plt.imshow(wordcloud, interpolation='bilinear')
7 plt.axis("off")
8 plt.show()
9
```

natural language processing nlp field artificial intelligence ai focuses interaction computers humans n  
atural language ultimate objective nlp enable computers understand interpret generate human language wa  
y meaningful useful.nlp techniques used analyze text extract relevant information derive insights unstr  
uctured data



**for text summarization**

```
In [20]: 1 # Tokenize the text into sentences
2 sentences = sent_tokenize(text)
3
4 # Tokenize each sentence into words and remove stopwords
5 stop_words = set(stopwords.words('english'))
6 word_frequencies = {}
7 for sentence in sentences:
8     word_tokens = word_tokenize(sentence)
9     filtered_words = [word.lower() for word in word_tokens if word.isalnum() and word.lower() not in
10     for word in filtered_words:
11         if word in word_frequencies:
12             word_frequencies[word] += 1
13         else:
14             word_frequencies[word] = 1
15 word_frequencies
```

```
Out[20]: {'natural': 2,
'language': 3,
'processing': 1,
'nlp': 2,
'field': 1,
'artificial': 1,
'intelligence': 1,
'ai': 1,
'focuses': 1,
'interaction': 1,
'computers': 2,
'humans': 1,
'ultimate': 1,
'objective': 1,
'enable': 1,
'understand': 1,
'interpret': 1,
'generate': 1,
'human': 1,
'way': 1,
'meaningful': 1,
'techniques': 1,
'used': 1,
'analyze': 1,
'text': 1,
'extract': 1,
'relevant': 1,
'information': 1,
'derive': 1,
'insights': 1,
'unstructured': 1,
'data': 1}
```

```

In [27]: 1 import nltk
2 from nltk.corpus import stopwords
3 from nltk.tokenize import word_tokenize, sent_tokenize
4
5 # Sample text
6 text = """
7 Natural language processing (NLP) is a field of artificial intelligence (AI) that focuses on the int
8 computers and humans through natural language.The ultimate objective of NLP is to enable computers t
9 interpret, and generate human language in a way that is both meaningful and useful.NLP techniques ar
10 text, extract relevant information, and derive insights from unstructured data.
11 """
12
13 # Tokenize the text into sentences
14 sentences = sent_tokenize(text)
15
16 # Tokenize each sentence into words and remove stopwords
17 stop_words = set(stopwords.words('english'))
18 word_frequencies = {}
19 for sentence in sentences:
20     word_tokens = word_tokenize(sentence)
21     filtered_words = [word.lower() for word in word_tokens if word.isalnum() and word.lower() not in
22     for word in filtered_words:
23         if word in word_frequencies:
24             word_frequencies[word] += 1
25         else:
26             word_frequencies[word] = 1
27
28 # Calculate sentence scores based on word frequencies
29 sentence_scores = {}
30 for sentence in sentences:
31     for word in word_tokenize(sentence.lower()):
32         if word in word_frequencies:
33             if len(sentence.split(' ')) < 30:
34                 if sentence in sentence_scores:
35                     sentence_scores[sentence] += word_frequencies[word]
36                 else:
37                     sentence_scores[sentence] = word_frequencies[word]
38
39 # Get top N sentences with highest scores
40 # Get top N sentences with highest scores
41 num_sentences = 2
42 summary_sentences = sorted(sentence_scores, key=sentence_scores.get, reverse=True)[:num_sentences]
43 summary = ' '.join(summary_sentences)
44 print("Summary: ", summary)
45
46 print("\noriginal text leghth:",len(text))
47 print("summarize text length:", len(summary))

```

Summary:

original text leghth: 433  
summarize text length: 0

## Set A: 4 (whatsapp chat analysis)

```
In [30]: 1 # reading text file
2 import re
3 import nltk
4 with open(r"C:\Users\hp\Desktop\TYCS\DATA ANALYTICS\_chat.txt", encoding="utf8") as chat:
5     chat_text1 = chat.read()
6
7 ### Cleaning and Preparing Data
8
9 # removing punct
10 from nltk.tokenize import word_tokenize
11 from string import punctuation
12
13 #removing all elements other than alphabets
14 chat_text=re.sub('[^a-zA-Z]', ' ', chat_text1)
15 # chat_text
16
17 #word tokenization
18 words = word_tokenize(chat_text)
19
20 #removing all stopwords
21 from nltk.corpus import stopwords,re
22 sw = set(stopwords.words("english"))
23 filterd_words = [w.lower() for w in words if w not in sw]
24
25 # removing names and other anomalies like am,pm
26 unwanted = {'vishnu','shahain','raghavan','preeti','pankaj','sinha','sharma',
27             'sahil','phatania','neha','mukti','omitted','image','kushbhu','am','pm'}
28 filterd_words = [ele for ele in filterd_words if ele not in unwanted]
29 # filterd_words
30
31 #wordcloud
32 import matplotlib.pyplot as plt
33 from wordcloud import WordCloud
34 text = ' '.join(filterd_words)
35 processed_text=re.sub('[^a-zA-Z]', ' ', text)
36 wordcloud = WordCloud().generate(processed_text)
37 # plt.figure(figsize=(10, 5))
38 plt.imshow(wordcloud, interpolation='bilinear')
39 plt.show()
40
41 # showing frequency distribution
42 from nltk import FreqDist
43 frequency_distribution= FreqDist(filterd_words)
44 frequency_distribution.most_common()
45
46 #sentiment analysis
47 from nltk.sentiment.vader import SentimentIntensityAnalyzer
48 vader_analyzer=SentimentIntensityAnalyzer()
49 print(vader_analyzer.polarity_scores( chat_text))
50
51 import matplotlib.pyplot as plt
52 frequency_distribution.plot(22,cumulative=False)
53 plt.show()
```



## Q1

```
In [61]: 1 import pandas as pd
2 import matplotlib.pyplot as plt
3 from wordcloud import WordCloud
4 from collections import Counter
5 from nltk.corpus import stopwords
6 from nltk.tokenize import word_tokenize
7
8 # Read the dataset
9 df = pd.read_csv(r"C:\Users\hp\Desktop\TYCS\DATA ANALYTICS\instagram_global_top_1000.csv")
10 df
```

Out[61]:

	Country	Rank	Account	Title	Link	Category	Followers
0	All	1	cristiano	Cristiano Ronaldo	<a href="https://www.instagram.com/cristiano/">https://www.instagram.com/cristiano/</a>	Sports with a ball	400100000.0
1	All	2	kyliejenner	Kylie ♥	<a href="https://www.instagram.com/kyliejenner/">https://www.instagram.com/kyliejenner/</a>	Fashion Modeling Beauty	308800000.0
2	All	3	leomessi	Leo Messi	<a href="https://www.instagram.com/leomessi/">https://www.instagram.com/leomessi/</a>	Sports with a ball Family	306300000.0
3	All	4	kendalljenner	Kendall	<a href="https://www.instagram.com/kendalljenner/">https://www.instagram.com/kendalljenner/</a>	Modeling Fashion	217800000.0
4	All	5	selenagomez	Selena Gomez	<a href="https://www.instagram.com/selenagomez/">https://www.instagram.com/selenagomez/</a>	Music Lifestyle	295800000.0
...	...	...	...	...	...	...	...
995	All	996	senoritasaeva	Dina Saeva	<a href="https://www.instagram.com/senoritasaeva/">https://www.instagram.com/senoritasaeva/</a>	Lifestyle Music Modeling	7700000.0
996	All	997	manuelneuer	Manuel Neuer	<a href="https://www.instagram.com/manuelneuer/">https://www.instagram.com/manuelneuer/</a>	Sports with a ball	11500000.0
997	All	998	sahilkhan	India's Youth & Fitness ICON@	<a href="https://www.instagram.com/sahilkhan/">https://www.instagram.com/sahilkhan/</a>	Fitness Gym	10100000.0
998	All	999	mohanshakti	Shakti Mohan	<a href="https://www.instagram.com/mohanshakti/">https://www.instagram.com/mohanshakti/</a>	Art Artists Cinema Actors/actresses	13700000.0
999	All	1000	eduincaz	Eduin Caz	<a href="https://www.instagram.com/eduincaz/">https://www.instagram.com/eduincaz/</a>	Lifestyle	6200000.0

1000 rows × 11 columns



```
In [62]: 1 # i. Find top 5 Instagram influencers from India
2 top_5_influencers_india = df[df['Audience Country'] == 'India'].nlargest(5, 'Followers')
3 print("Top 5 Instagram influencers from India:")
4 sorted_top_5_influencers_india = top_5_influencers_india.sort_values(by='Followers', ascending=False)
5
6 print(sorted_top_5_influencers_india[['Account', 'Followers']])
```

Top 5 Instagram influencers from India:

```
Account Followers
28 instagram 469600000.0
0 cristiano 400100000.0
27 therock 295800000.0
19 justinbieber 219800000.0
14 virat.kohli 182600000.0
```

```
In [63]: 1 # ii. Find Instagram account with least number of Followers
2 least_Followers_account = df.nsmallest(1, 'Followers')
3 print("\nInstagram account with least number of Followers:")
4 print(least_Followers_account[['Account', 'Followers']])
```

```
Instagram account with least number of Followers:
Account Followers
747 yooney1 2800000.0
```

```
In [64]: 1 # Tokenize each string in the 'Category' column
2 categories = df['Category'].astype(str)
```

```
In [65]: 1 df["Category"]
```

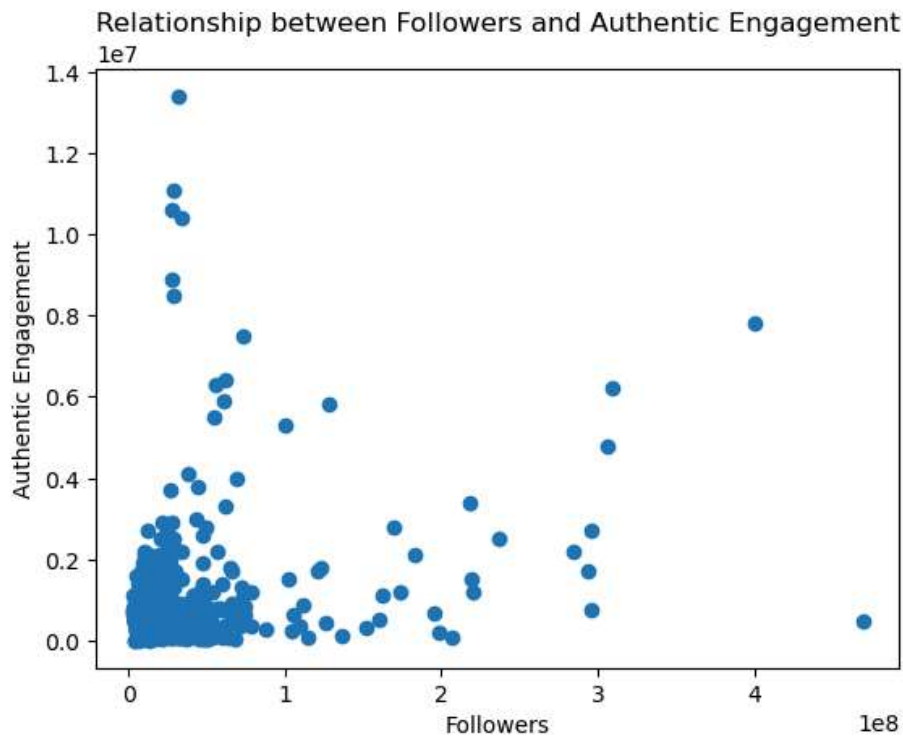
```
Out[65]: 0 Sports with a ball
1 Fashion|Modeling|Beauty
2 Sports with a ball|Family
3 Modeling|Fashion
4 Music|Lifestyle
...
995 Lifestyle|Music|Modeling
996 Sports with a ball
997 Fitness|Gym
998 Art|Artists|Cinema|Actors/actresses
999 Lifestyle
Name: Category, Length: 1000, dtype: object
```

```
In [66]: 1 # Tokenize each string in the 'Category' column
2 categories = df['Category'].astype(str)
3 # Assuming 'categories' is the column you want to tokenize
4 # Tokenize each string in the 'categories' column
5 word_tokens = categories.apply(lambda x: word_tokenize(str(x)))
6 word_tokens
```

```
Out[66]: 0 [Sports, with, a, ball]
1 [Fashion|Modeling|Beauty]
2 [Sports, with, a, ball|Family]
3 [Modeling|Fashion]
4 [Music|Lifestyle]
...
995 [Lifestyle|Music|Modeling]
996 [Sports, with, a, ball]
997 [Fitness|Gym]
998 [Art|Artists|Cinema|Actors/actresses]
999 [Lifestyle]
Name: Category, Length: 1000, dtype: object
```



```
In [70]: 1 # v. Visualize relationship between Followers and Authentic engagement
2 plt.scatter(df['Followers'], df['Authentic engagement'])
3 plt.title('Relationship between Followers and Authentic Engagement')
4 plt.xlabel('Followers')
5 plt.ylabel('Authentic Engagement')
6 plt.show()
7
```



```
In [71]: 1 correlation = df['Followers'].corr(df['Authentic engagement'])
2 correlation
```

Out[71]: 0.2870667818696664

## Q2

```
In [53]: 1 import pandas as pd
2 from nltk.tokenize import word_tokenize
3 from nltk.sentiment.vader import SentimentIntensityAnalyzer
4 import nltk
5 # nltk.download('vader_lexicon')
6
7 # Load the dataset
8 df = pd.read_csv(r"C:\Users\hp\Desktop\TYCS\DATA ANALYTICS\covid_2021_1.csv")
9 print(df.head())
```

```

                                query \
0 coronavirus|covid%7C19|pandemic|vaccine
1 coronavirus|covid%7C19|pandemic|vaccine
2 coronavirus|covid%7C19|pandemic|vaccine
3 coronavirus|covid%7C19|pandemic|vaccine
4 coronavirus|covid%7C19|pandemic|vaccine

                                url \
0 https://www.youtube.com/watch?v=LyRoIKjmk9o (https://www.youtube.com/watch?v=LyRoIKjmk9o)
1 https://www.youtube.com/watch?v=LyRoIKjmk9o (https://www.youtube.com/watch?v=LyRoIKjmk9o)
2 https://www.youtube.com/watch?v=LyRoIKjmk9o (https://www.youtube.com/watch?v=LyRoIKjmk9o)
3 https://www.youtube.com/watch?v=LyRoIKjmk9o (https://www.youtube.com/watch?v=LyRoIKjmk9o)
4 https://www.youtube.com/watch?v=dwMWM1F7NMI (https://www.youtube.com/watch?v=dwMWM1F7NMI)

                                title                upload_date \
0 Which Coronavirus Vaccine Is More Effective? C... 2021-01-01T09:04:16Z
1 Which Coronavirus Vaccine Is More Effective? C... 2021-01-01T09:04:16Z
2 Which Coronavirus Vaccine Is More Effective? C... 2021-01-01T09:04:16Z
3 Which Coronavirus Vaccine Is More Effective? C... 2021-01-01T09:04:16Z
4 Growing Number Of Americans Going Hungry Durin... 2021-01-01T11:59:49Z

                                channel  views  likes  dislikes  comment_count \
0 India Today  8276    58     7           5
1 India Today  8276    58     7           5
2 India Today  8276    58     7           5
3 India Today  8276    58     7           5
4 NBC News    74546   874    62          601

                                comment_text \
0
1                                OMG 🍷🍷🍷🍷🍷🍷
2                                I love my Indian Mumbai me welcome
3                                What about indian vaccine?
4                                2:32 I like that 🍷🍷🍷🍷🍷🍷
5                                And all the rich companies that have profited ...

                                comment_author  comment_date  comment_likes \
0                                Brendan Eric  2021-01-01T09:23:59Z  0
1 Bhallesingh Rathore rathore  2021-01-01T09:06:51Z  1
2 Kudip singh hong kong.  2021-02-03T08:17:14Z  0
3 Bo Brice  2021-01-01T09:55:34Z  0
4 Poppy Jalto  2021-01-01T13:17:56Z  93

                                DATE
0 2021-01-01 09:04:16+00:00
1 2021-01-01 09:04:16+00:00
2 2021-01-01 09:04:16+00:00
3 2021-01-01 09:04:16+00:00
4 2021-01-01 11:59:49+00:00
```

```
In [54]: 1 # i. Data Cleaning
2 # Drop rows with missing values
3 df.dropna(inplace=True)
```

```
In [55]: 1 df['tokenized_comments'] = df['comment_text'].apply(word_tokenize)
2 df['tokenized_comments']
```

```
Out[55]: 0 [OMG, 🤪👄, 🍷❤️❤️]
1 [I, love, my, Indian, Mumbai, me, welcome]
2 [What, about, indian, vaccine, ?]
3 [2:32, I, like, that 🤪👄, 🍷❤️❤️]
4 [And, all, the, rich, companies, that, have, p...
...
41583 [Thanks, da]
41584 [👍👍, very, helpful, information]
41585 [👍]
41586 [Very, very, much, informative, vlog, .., Sunn...
41587 [Thanking, you, for, this, information]
Name: tokenized_comments, Length: 41588, dtype: object
```

```
In [56]: 1 # iii. Sentiment Analysis
2 def analyze_sentiment(comment):
3     analyzer = SentimentIntensityAnalyzer()
4     sentiment_scores = analyzer.polarity_scores(comment)
5     if sentiment_scores['compound'] > 0:
6         return 'positive'
7     elif sentiment_scores['compound'] < 0:
8         return 'negative'
9     else:
10        return 'neutral'
11
12 df['sentiment'] = df['comment_text'].apply(analyze_sentiment)
13 df['sentiment']
```

```
Out[56]: 0 neutral
1 positive
2 neutral
3 positive
4 positive
...
41583 positive
41584 positive
41585 neutral
41586 positive
41587 neutral
Name: sentiment, Length: 41588, dtype: object
```

```
In [ ]: 1 # Calculate percentage of positive, negative, and neutral comments
2 sentiment_counts = df['sentiment'].value_counts(normalize=True) * 100
3
4 print("Percentage of Comments by Sentiment:")
5 print(sentiment_counts)
6
```

**Q3**

```
In [73]: 1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 # i. Read the dataset and perform data cleaning operations
5 df = pd.read_csv('INvideos.csv')
6
7 # Drop duplicates and missing values
8 df.drop_duplicates(inplace=True)
9 df.dropna(inplace=True)
10
11
```

```
In [74]: 1 # ii. Find total views, likes, dislikes, and comment count
2 total_views = df['views'].sum()
3 total_likes = df['likes'].sum()
4 total_dislikes = df['dislikes'].sum()
5 total_comments = df['comment_count'].sum()
6
7 print("Total Views:", total_views)
8 print("Total Likes:", total_likes)
9 print("Total Dislikes:", total_dislikes)
10 print("Total Comments:", total_comments)
```

Total Views: 32792912293  
Total Likes: 845128501  
Total Dislikes: 51994375  
Total Comments: 83413188

```
In [75]: 1 # iii. Find the Least and topmost Liked and commented videos
2
3 #row corresponding to the Least number of likes
4 least_liked_video = df.loc[df['likes'].idxmin()]
5 #row corresponding to the max number of Likes
6 top_liked_video = df.loc[df['likes'].idxmax()]
7
8 #row corresponding to the Least number of comments
9 least_commented_video = df.loc[df['comment_count'].idxmin()]
10 #row corresponding to the max comments
11 top_commented_video = df.loc[df['comment_count'].idxmax()]
```

```
In [76]: 1 print("\nLeast Liked Video:\n",least_liked_video[['title', 'likes']])
2 print("\nTop Liked Video:\n",top_liked_video[['title', 'likes']])
3 print("\nLeast Commented Video:\n",least_commented_video[['title', 'comment_count']])
4 print("\nTop Commented Video:\n",top_commented_video[['title', 'comment_count']])
```

Least Liked Video:

```
title    Breaking News IT Raid - நடந்தது என்ன? சிக்கியத...
likes                                         0
Name: 34, dtype: object
```

Top Liked Video:

```
title    YouTube Rewind: The Shape of 2017 | #YouTubeRe...
likes                                         2912710
Name: 5408, dtype: object
```

Least Commented Video:

```
title    पीरियड्स के समय, पेट पर पति करता ऐसा, देखकर दं...
comment_count                                0
Name: 1, dtype: object
```

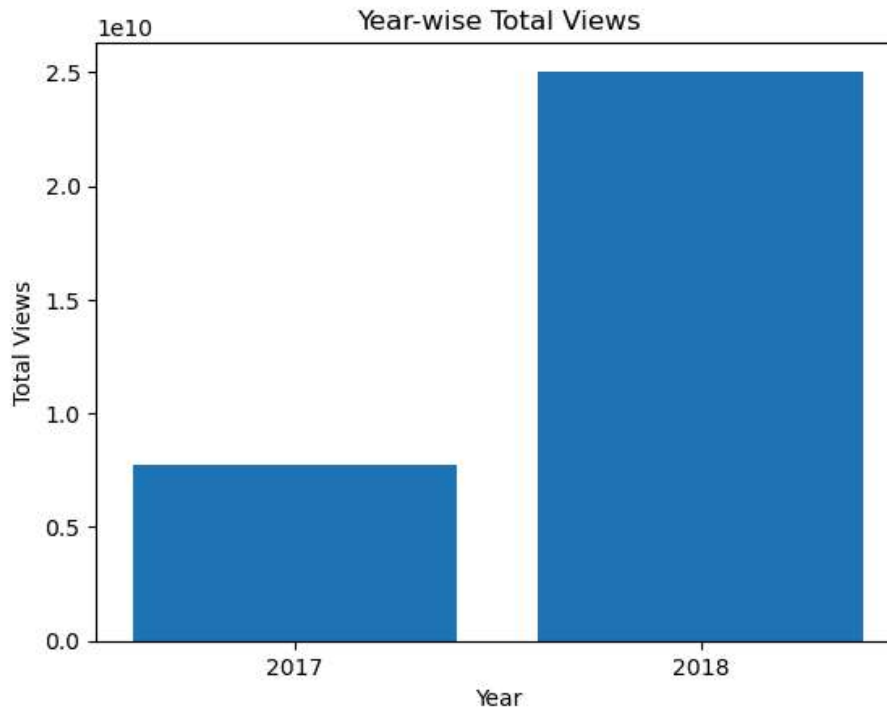
Top Commented Video:

```
title    YouTube Rewind: The Shape of 2017 | #YouTubeRe...
comment_count                                827755
Name: 4743, dtype: object
```

```
In [79]: 1 # iv. Perform year-wise statistics for views and plot the analyzed data
2 df['trending_date'] = pd.to_datetime(df['trending_date'], format='%y.%d.%m')
3 df['year'] = df['trending_date'].dt.year
4 yearly_views = df.groupby('year')['views'].sum()
5 df['trending_date']
```

```
Out[79]: 0      2017-11-14
1      2017-11-14
2      2017-11-14
3      2017-11-14
4      2017-11-14
...
37300  2018-06-14
37301  2018-06-14
37302  2018-06-14
37319  2018-06-14
37330  2018-06-14
Name: trending_date, Length: 32562, dtype: datetime64[ns]
```

```
In [80]: 1 plt.bar(yearly_views.index, yearly_views.values)
2 plt.xlabel('Year')
3 plt.ylabel('Total Views')
4 plt.title('Year-wise Total Views')
5 plt.xticks(yearly_views.index)
6 plt.show()
```



```
In [ ]: 1 # v. Plot the viewers who reacted on videos
2 reacted_viewers = df[['likes', 'dislikes']]
3 reacted_viewers.plot(kind='bar', stacked=True)
4 plt.xlabel('Video Index')
5 plt.ylabel('Number of Viewers')
6 plt.title('Viewers who Reacted on Videos')
7 plt.show()
8
```

```
In [ ]: 1
```