Assignment 3 **Text and Social Media Analytics** Set A Q1. Consider any text paragraph preprocess the text to remove any special characters and digits. Generate the summary using extractive

summarization process.

In []:
!pip install nltk import nltk

]:
nltk.download('all')

In []:
 nltk.download('punkt') nltk.download('stopwords')

In [12]:
text = """ Artificial intelligence is the future. Artificial intelligence is science fiction. Artificial intell

In [30]:
import nltk from nltk.corpus import stopwords from nltk.tokenize import word_tokenize, sent_tokenize

In [31]:
import re import heapq def clean data(data): text = $re.sub(r"\[[0-9]*\]"," ", data)$ text = text.lower()

 $text = re.sub(r'\s+'," ", text)$ text = re.sub(r", ", " ", text)return text cleaned data = clean data(text) sent tokens = sent tokenize(cleaned data) word tokens = word tokenize(cleaned_data) word frequency = {} stopwords = set(stopwords.words("english"))

for word in word tokens: if word not in stopwords: if word not in word frequency.keys(): word frequency[word]=1 word frequency[word] +=1

else: maximum frequency = max(word frequency.values()) for word in word frequency.keys(): word_frequency[word] = (word_frequency[word]/maximum_frequency) sentences score = {} for sentence in sent tokens: for word in word tokenize(sentence): if word in word_frequency.keys():

if (len(sentence.split(" "))) <30:</pre> if sentence not in sentences score.keys(): sentences score[sentence] = word frequency[word] sentences score[sentence] += word frequency[word] def get key(val): for key, value in sentences_score.items(): if val == value:

key = get_key(max(sentences_score.values())) summary = heapq.nlargest(5, sentences score, key=sentences score.get) print(" ".join(summary)) artificial intelligence is already part of our everyday lives. and all three are part of the reason why alphago trounced lee se-dol. all those statements are true it just depends on what flavor of ai you are referring to. artificial intelligence is science fiction. artificial intelligence is the future. Set A Q2. Consider any text paragraph remove the stopwords. Tokenize the paragraph to extract words and sentences. Calculate the word frequency distribution and plot the frequency. Plot the wordcloud of the text. In [51]:
import nltk from nltk.corpus import stopwords from nltk.tokenize import sent_tokenize, word_tokenize

from nltk.probability import FreqDist In [52]:
text="""As the process of analyzing raw data to find trends and answer questions, the definition of data analyt In [53]:
 sent_tokens = sent_tokenize(text) word tokens = word tokenize(text) stopwords = set(stopwords.words("english")) frequency distribution=FreqDist(word tokens) print(frequency distribution) frequency distribution <FreqDist with 30 samples and 38 outcomes> Out[53]:

'many': 2, 'As': 1, 'process': 1, 'analyzing': 1, 'ra w':

1, ...})

FreqDist({'the': 3, 'of': 3, 'data': 2, ',': 2, '.': 2, In [54]:
 import matplotlib.pyplot as plt ax = plt.axes() frequency distribution.plot(20,cumulative=False) #first 20 words starting from maximum onn Xaxis ax.set title('Frequency Plot of words') plt.show() 3.00 2.75 2.50 2.25 2.00 1.75 1.50 1.25

In []:
!pip install wordcloud from wordcloud import WordCloud In [74]:
 word cloud=WordCloud(background_color='black').generate(text) In [75]:
plt.figure() plt.imshow(word cloud, interpolation='bilinear') plt.show()

Samples

print("Overall rating for sentence is Positive")

print("Overall rating for sentence is Negative")

print("Overall rating for sentence is neutral")

The sentence is rated as 0.0 % Negative The sentence is rated as 71.0 % Neutral Overall rating for sentence is Positive

sent tokens = sent tokenize(text) word tokens = word tokenize(text)

stopwords = set(stopwords.words("english"))

'myself', 'shan', 'off', 'her', 'now'}

filtered words list.append(words)

'many', 'different', 'goals', '.']

from nltk.stem import PorterStemmer

for words in filtered words list:

from nltk.stem import PorterStemmer

PorterStemmer().stem('eaten')

porter stemmer=PorterStemmer()

stemmed_text_words=[]

Filtered Words:

'goal', '.']

print("Tokenized Words : \n", word tokens, "\n")

print("Filtered Words : \n", filtered_words_list,"\n")

stemmed_text_words.append(porter_stemmer.stem(words))

print("Stemmed Words : \n", stemmed text words, "\n")

print("Filtered Words : \n", word_tokens,"\n")

for words in word tokens: if words not in stopwords:

Tokenized Words:

Filtered Words :

Tokenized Sentences :

'the', '

Tokenized Words :

print(stopwords)

The sentence is rated as 28.999999999999 % Positive

print("Tokenized Sentences: \n", sent tokens, "\n") #break paragraph into sentences (. full stop is considere

['As the process of analyzing raw data to find trends and answer questions, the definition of data analytics c

s', ',', 'the', 'definition', 'of', 'data', 'analytics', 'captures', 'its', 'broad', 'scope', 'of', 'the', 'fie ld', '.', 'However', ',', 'it', 'includes', 'many', 'techniques', 'with', 'many', 'different', 'goals', '.']

{'or', 'about', 'above', 'not', 'o', 'than', 'wasn', 'yourselves', 'down', 'most', "aren't", 'mustn', 'too', 's ame', 'such', 'weren', 'these', 'from', 'very', 'why', "shan't", 'on', 'll', "hadn't", 'over', 'd', 'she', 'bel ow', 'them', 'through', 've', 'herself', 'don', 'for', 'be', 'until', 'a', "haven't", "wouldn't", 'they', 'thei rs', 'i', 'when', 'hadn', 'ma', 'few', 'by', 'hers', 'couldn', 'it', 'under', 'no', "she's", 'because', 'thei r', 'should', 'just', 'after', 'he', 'what', 'will', 'needn', 'do', "it's", "won't", 'have', 'didn', 'did', 't', 'won', 'itself', "doesn't", 'while', 'himself', 're', 'nor', 'is', 'in', "isn't", 'once', 'so', 'ours', 'y ourself', 'hasn', 'only', 'yours', 'ain', "wasn't", 'shouldn', "you'll", 'y', 'out', 'm', 'to', "couldn't", "yo u'd", 'are', 'doesn', 'between', 'any', 'aren', 'been', 'haven', 'my', 'we', 'here', 'who', "don't", "that'll", 'which', 'more', 'into', 'if', "weren't", 'and', 'up', "shouldn't", 'again', 'of', 'against', 'him', 'other', 'ourselves', 'me', "you're", 'those', 'were', 'during', "you've", 'but', 'an', 'where', 'was', 'then', 'does', "mightn't", 's', 'this', "should've", 'your', 'before', 'both', 'had', 'the', 'further', 'having', 'can', 'is n', 'each', 'wouldn', 'with', "needn't", 'his', 'has', 'how', "didn't", 'being', 'as', "hasn't", 'mightn', 'ou r', "mustn't", 'am', 'themselves', 'some', 'doing', 'at', 'own', 'you', 'whom', 'there', 'all', 'that', 'its',

['As', 'the', 'process', 'of', 'analyzing', 'raw', 'data', 'to', 'find', 'trends', 'and', 'answer', 'question s', ',', 'the', 'definition', 'of', 'data', 'analytics', 'captures', 'its', 'broad', 'scope', 'of', 'the', 'fie ld', '.', 'However', ',', 'it', 'includes', 'many', 'techniques', 'with', 'many', 'different', 'goals', '.']

['As', 'process', 'analyzing', 'raw', 'data', 'find', 'trends', 'answer', 'questions', ',', 'definition', 'dat a', 'analytics', 'captures', 'broad', 'scope', 'field', '.', 'However', ',', 'includes', 'many', 'techniques',

['As', 'the', 'process', 'of', 'analyzing', 'raw', 'data', 'to', 'find', 'trends', 'and', 'answer', 'question s', ',', 'the', 'definition', 'of', 'data', 'analytics', 'captures', 'its', 'broad', 'scope', 'of', 'the', 'fie ld', '.', 'However', ',', 'it', 'includes', 'many', 'techniques', 'with', 'many', 'different', 'goals', '.']

['as', 'process', 'analyz', 'raw', 'data', 'find', 'trend', 'answer', 'question', ',', 'definit', 'data', 'ana lyt', 'captur', 'broad', 'scope', 'field', '.', 'howev', ',', 'includ', 'mani', 'techniqu', 'mani', 'differ',

'process', 'of', 'analyzing', 'raw', 'data', 'to', 'find', 'trends', 'and', 'answer', 'question

print("Tokenized Words: \n", word tokens, "\n") #break sentences into words (<spaces are considered >)

aptures its broad scope of the field.', 'However, it includes many techniques with many different goals.']

Set A Q3. Consier the following review message .perform sentiment analysis on the message In []:
nltk.download('vader_lexicon') In [77]:
 from nltk.sentiment.vader import SentimentIntensityAnalyzer vader analyzer=SentimentIntensityAnalyzer() In [78]:
 text1="I purchased headphones online. I am very happy with the product." print(vader analyzer.polarity scores(text1))

{'neg': 0.0, 'neu': 0.667, 'pos': 0.333, 'compound': 0.6115} In [79]:
 text2="I saw the movie yesterday. The animation was really good but the script was ok." print(vader analyzer.polarity scores(text2)) {'neg': 0.0, 'neu': 0.71, 'pos': 0.29, 'compound': 0.5989} In [81]:
 result2=vader_analyzer.polarity_scores(text2) # To find percentage of ratings print("The sentence is rated as ",result1['pos']*100,"% Positive") print("The sentence is rated as ",result1['neg']*100,"% Negative") print("The sentence is rated as ",result1['neu']*100,"% Neutral") if result1['compound']>=0.05: elif result1['compound'] <=-0.05:</pre> else:

In [90]:
 text="""As the process of analyzing raw data to find trends and answer questions, the definition of data analyt In [92]:
 from nltk.tokenize import sent_tokenize, word_tokenize

In [95]:
 from nltk.corpus import stopwords

In [97]:
filtered_words_list=[] In [98]:
#Stemming change writing, wrote, written can stemmed or reduced as write In [103" #Stemming

'eaten' Out[103 In [104"PorterStemmer().stem('eating') 'eat' Out[104 In [107 PorterStemmer().stem('eated') Q. Find the lemmatized Out[107... word for eaten. 'eat' In []:
#Lemmatization import nltk nltk.download('wordnet')

In [106"# Lemmatization

Lemmatization
from nltk.stem.wordnet import WordNetLemmatizer
lemmatizer=WordNetLemmatizer()
word_text="eaten"
print("Lemmatized Word :

",lemmatizer.lemmatize(word_text,"v"))Lemmatized Word : eat