,	Tasks carried out: Data Extraction Data Exploration and EDA. Text Processing. Model Building Creating a pickle file for model deployment using web application:
•	1) Importing The Required Libraries: import numpy as np import pandas as pd import matplotlib.pyplot as plt %matplotlib inline import datetime #import nltk import string
	<pre>#nltk.download('stopwords') from nltk.corpus import stopwords, wordnet from nltk.tokenize import word_tokenize from nltk.stem import WordNetLemmatizer #from gensim.models import Word2Vec #from gensim.models import KeyedVectors from sklearn.feature_extraction.text import CountVectorizer from sklearn.feature_extraction.text import TfidfVectorizer from sklearn.metrics import pairwise_distances from sklearn.metrics.pairwise import euclidean_distances #from sklearn.metrics.pairwise import cosine_similarity #from scipy.sparse import csr_matrix</pre>
	<pre>import pickle import warnings warnings.filterwarnings("ignore") 2) Loading The Dataset: news_articles = pd.read_csv(r"C:\Users\sksho\Desktop\ZenteiQ\Work\Data\News category dataset.csv") news_articles.head() category</pre>
	CRIME There Were 2 Mass Shootings In Texas Last Week Melissa Jeltsen https://www.huffingtonpost.com/entry/texas-ama She left her husband. He killed their children 2018-05-26 1 ENTERTAINMENT Will Smith Joins Diplo And Nicky Jam For The 2 Andy McDonald https://www.huffingtonpost.com/entry/will-smit Of course it has a song. 2 ENTERTAINMENT Hugh Grant Marries For The First Time At Age 57 Ron Dicker https://www.huffingtonpost.com/entry/hugh-gran The actor and his longtime girlfriend Anna Ebe 2018-05-26 3 ENTERTAINMENT Jim Carrey Blasts 'Castrato' Adam Schiff And D Ron Dicker https://www.huffingtonpost.com/entry/jim-carre The actor gives Dems an ass-kicking for not fi 2018-05-26 4 ENTERTAINMENT Julianna Margulies Uses Donald Trump Poop Bags Ron Dicker https://www.huffingtonpost.com/entry/julianna The "Dietland" actress said using the bags is 2018-05-26 1 The "Dietland" actress said using the bags is 2018-05-26
	Colass 'pandas.core.frame.DataFrame'> RangeIndex: 200853 entries, 0 to 200852 Data columns (total 6 columns):
	Note: From above we can observe that there are some missing values. Hence we will check the Number of missing values in the data and the contribution of missing values in the data. missing_values_contribution = pd.DataFrame(('No.of Missing Values': news_articles.isna().sum(),
	category 0 0.000000 headline 6 0.002987 authors 36620 18.232239 link 0 0.000000 short_description 19712 9.814143 date 0 0.000000
	Note: There are total 200853 records in the data with 6 object type features. Also there are missing values in the data. • We can observe that the feature authors have maximum number of missing values followed by shrot_description and headlines which contributes about 18.23%, 9.81%, and 0.0029% of data respectively. • We will drop this null values as this will not contribute much in our analysis of dataset. news_articles.dropna (inplace=True) news_articles.isna() .sum() category 0 headline 0
	authors 0 link 0 short_description 0 date 0 dtype: int64 news_articles.info() <class 'pandas.core.frame.dataframe'=""> Int64Index: 148983 entries, 0 to 200848 Bata columns (total 6 columns): # Column Non-Null Count Dtype</class>
j	O category 148983 non-null object 1 headline 148983 non-null object 2 authors 148983 non-null object 3 link 148983 non-null object 4 short_description 148983 non-null object 5 date 148983 non-null object dtypes: object(6) memory usage: 8.0+ MB Note: we can observe that all null values have be dropped and we are left with 148983 non null values.
	news_articles.shape (148983, 6) Note: From the dataset we can note that the columns category, authors, and short description can have same values. hence we have to note that if there are any duplicate values in the headlines column. and eliminate them. news_articles['headline'].count() 148983
	Takegry Salas Shootings In Texas Last Week Melissa Jeltse Naty Melissa Melissa Melissa Melissa Melissa Jeltse Naty Melissa Melissa Melissa Melissa Jeltse Naty Melissa
	in.ii.ii.ii.iii.200843TECHGood Games Is It possible?Mateo Gutierrez, Contributon/nArtisishttps://www.huffingtonpost.com/entry/google-pl
	news_articles.sort_values('headline', inplace=True, ascending=False) duplicated_articles_series = news_articles.duplicated('headline', keep = False) news_articles = news_articles[~duplicated_articles_series] print("Total number of articles after removing duplicates:", news_articles['headline'].nunique()) Total number of articles after removing duplicates: 147706 To see the total number of unique records. news_articles.nunique()
	category 41 headline 147706 authors 27055 link 147706 short_description date 2309 dtype: int64 3) Exploratory Data Analysis:
	news_articles["category"].value_counts().head() POLITICS
	COLLEGE 859 ENVIRONMENT 669 CULTURE & ARTS 611 Name: category, dtype: int64 plt.figure(figsize=(10,6)) news_articles("category"].value_counts().plot.bar(xlabel = 'News Category', ylabel = 'Number of Articles', title = 'Number of Articles Per Category') plt.show()
	Number of Articles Per Category 25000 - 20000 -
	10000 - 10000
	WELLNESS ENTERTAINMENT TRAVEL STYLE & BEAUTY PARENTING HEALTHY LIVING QUEER VOICES FOOD & DRINK BLACK VOICES GOMEDY SPORTS PARENTS RELIGION WEIND NEWS RELIGION WEND NEWS RELIGION WEND NEWS FIFTY GOOD NEWS EDUCATION WENDOST FIFTY GOOD NEWS EDUCATION TECH STYLE WORLDPOST FIFTY GOOD NEWS EDUCATION TAKEN TAKEN TAKEN TAKEN TO THE WORLDPOST TAKEN TO
	Note: • We can observe that Top 5 articles that are watched are based on Politics, Wellness, Entertainment, Travel, Style & Beauty. • We can observe that Bottom 5 articles that are watched are based on Culture & Arts, Environment, College, Arts, Latino Voices. **Rews_articles['date'] = pd.to_datetime (news_articles['date'])
	# Resample the DataFrame by Year and count the number of articles news_articles_monthly = news_articles.resample('Y', on='date')['headline'].count() # Create a bar plot of the monthly article counts plt.figure(figsize = (10,6)) plt.barh(news_articles_monthly.index.strftime('%Y'), news_articles_monthly) # Set the plot title and axis labels plt.title('Number of Articles per Year') plt.xlabel('Year') plt.ylabel('Number of Articles')
	# Display the plot plt.show() Number of Articles per Year 2018 - 2017 -
	2017 - 2016 - 2015 - 2014 - 20
	2013 - 2012 - 0 5000 10000 15000 20000 25000 Year
	Note: From above ditribution we can observe the total number of articles on yearly basis. news_articles.index = range(news_articles.shape[0]) news_articles["day and month"] = news_articles["date"].dt.strftime("%a") + "_" + news_articles["date"].dt.strftime("%b") news_articles["user_id"] = ["user_" + str(i) for i in range(1, 147707)]
	news_articles = news_articles.reindex (columns=['user_id', 'category', 'headline', 'authors', 'link', 'short_description', 'date', 'day and month']) news_articles user_id
	WEIRD NEWS
1	147706 rows × 8 columns news_articles_temp = news_articles.copy() news_articles_temp.head(1) user_id category
	news_articles_temp.tail(1) user_id category headline authors link short_description date day and month 147705 user_147706 ENTERTAINMENT "2 Guns": Two Pros In Action Jackie K Cooper, Contributor\nFilm Critic https://www.huffingtonpost.com/entry/2-guns-tw It helps that the supporting cast is almost as 2013-08-03 Sat_Aug Note: We can observe that all the changes has been done.
	4) Text Processing: stop_words = set(stopwords.words('english')) for i, headline in enumerate(news_articles_temp["headline"]): cleaned_words = [word.lower() for word in headline.split() if word.isalnum() and word.lower() not in stop_words] cleaned_headline = ' '.join(cleaned_words) news_articles_temp.at[i, "clean_headline"] = cleaned_headline
	<pre>news_articles_temp.at[i, "clean_headline"] = cleaned_headline if(i%1000==0): print(i) # To track number of records processed 0 1000 2000 3000 4000 5000</pre>
	6000 7000 8000 9000 10000 11000 12000 13000 14000 15000 15000
	17000 18000 19000 20000 21000 22000 22000 23000 24000 25000
	26000 27000 28000 28000 29000 30000 31000 32000 32000 33000 34000 35000
	36000 37000 38000 39000 40000 41000 42000 43000 44000 44000 45000 45000 460000 47000
	4800 4900 5000 5100 52000 53000 54000 55000 55000 55000 55000 56000 56000 56000 56000
	59000 60000 61000 62000 63000 64000 65000 65000 66000 67000 68000 69000
	7000 7100 7200 73000 74000 75000 76000 77000 78000 78000 78000 78000 78000
	81000 82000 83000 84000 85000 85000 86000 87000 88000 89000 99000
	93000 94000 95000 95000 96000 97000 98000 99000 1000000 1010000 102000 103000
	104000 105000 106000 107000 108000 109000 110000 111000 112000 113000 114000
	115000 116000 117000 118000 119000 120000 121000 122000 123000 124000 125000
	126000 127000 128000 129000 130000 131000 132000 132000 133000 134000 135000 135000
	137000 138000 140000 141000 142000 144000 144000 145000 145000 146000
	<pre>lemmatizer = WordNetLemmatizer() for i, headline in enumerate(news_articles_temp["headline"]): if isinstance(headline, str): lemmatized_words = [lemmatizer.lemmatize(word, pos="v") for word in word_tokenize(headline)] lemmatized_headline = ' '.join(lemmatized_words) news_articles_temp.at[i, "lemmatized_headline"] = lemmatized_headline if i % 1000 ==0: print(i) 0 1000</pre>
	2000 3000 4000 5000 6000 7000 8000 9000 11000 11000
	13000 14000 15000 16000 17000 18000 18000 20000 21000 22000 23000 24000
	24000 25000 26000 27000 28000 29000 30000 31000 31000 32000 33000 34000 35000
	36000 37000 38000 39000 40000 41000 42000 43000 44000 44000 45000 45000 45000 47000
	48000 49000 51000 51000 52000 53000 54000 55000 55000 56000 57000 58000 59000
	59000 61000 62000 63000 64000 65000 65000 66000 67000 68000 68000 69000 70000
	71000 72000 73000 74000 75000 75000 76000 77000 78000 78000 80000 80000 81000
	82000 83000 84000 85000 86000 87000 88000 89000 99000 90000
	93000 94000 95000 96000 97000 98000 990000 1000000 1010000 102000 1030000 1040000
	105000 106000 108000 108000 109000 110000 1110000 1120000 1120000 1130000 1140000 1150000
	115000 117000 118000 119000 120000 121000 122000 122000 123000 124000 125000 125000 125000 125000
	128000 129000 130000 131000 132000 132000 133000 134000 135000 136000 137000 138000 138000 138000
	14000 142000 143000 144000 145000 147000 5) Model Building:
	5.1 Using Bag of Words: headline_vectorizer = CountVectorizer() headline_features = headline_vectorizer.fit_transform(news_articles_temp['headline']) headline_features.get_shape() (147706, 48689)
	<pre>from sklearn.metrics.pairwise import euclidean_distances euclidean_sim = None # Declare as global variable def bag_of_words_based_model(user_id, num_similar_items): global euclidean_sim # Declare as global to modify the variable # Calculate the Euclidean similarity euclidean_sim = euclidean_distances(headline_features, headline_features[user_id].reshape(1, -1)) indices = np.argsort(euclidean_sim.ravel())[0:num_similar_items]</pre>
	<pre>df = pd.DataFrame({ 'publish_date': news_articles['date'][indices].values, 'category': news_articles['category'][indices].values, 'headline': news_articles['headline'][indices].values, 'link': news_articles['link'][indices].values, 'short_description': news_articles['short_description'][indices].values, 'Euclidean similarity with the queried article': euclidean_sim[indices].ravel() }) print("*" * 30, "Queried article details", "*" * 30) print('headline: ', news_articles['headline'][indices[0]])</pre>
	<pre>bag_of_words_based_model(1, 5) ***********************************</pre>

	Wellnes with Bone of Did to the
	Parameter Para
	Publish date State
	Decided the set of the standard of the standard for the s
	Part